

UNCLASSIFIED

**Project Report
LSP-156**

**Low-Power Embedded Analytics:
FY15 Line-Supported Information,
Computation, and Exploitation Program**

H.T. Nguyen
J.B. Muldavin
V.N. Gadepally
Arvind

TBD November 2015

Lincoln Laboratory
MASSACHUSETTS INSTITUTE OF TECHNOLOGY
LEXINGTON, MASSACHUSETTS



Prepared for the Assistant Secretary of Defense for Research and Engineering
under Air Force Contract FA8721-05-C-0002.

UNCLASSIFIED

This report is based on studies performed at Lincoln Laboratory, a federally funded research and development center operated by Massachusetts Institute of Technology. This work is sponsored by the Assistant Secretary of Defense for Research and Engineering under Air Force Contract FA8721-05-C-0002.

UNCLASSIFIED

**Massachusetts Institute of Technology
Lincoln Laboratory**

**Low-Power Embedded Analytics:
FY15 Line-Supported Information, Computation, and Exploitation Program**

*H.T. Nguyen
J.B. Muldavin
Group 102*

*V.N. Gadepally
Group 53*

*Prof. Arvind
MIT Computer Science and Artificial Intelligence Laboratory*

Project Report LSP-156

TBD November 2015

Lexington

Massachusetts

UNCLASSIFIED

This page intentionally left blank.

EXECUTIVE SUMMARY

This report covers the second year of the low-power embedded analytics project, a three-year university collaboration between Lincoln Laboratory and Professor Arvind's group at the MIT Computer Science and AI Laboratory (CSAIL). The goal of the project is to design and prototype a novel architecture that has wide potential applicability to important applications ranging from back-office big-data analytics to fieldable hot-spot systems providing storage-processing-communication services for off-grid sensors. Speed and power efficiency are the key metrics.

Current state-of-the-art approaches for big-data aim toward scaling out to many computers to meet processing, storage capacity, and access bandwidth requirements. Data is distributed over many computers, and complex processing is decomposed into tasks that operate on localized data and aggregated back together. With an emphasis on scalability and cross-platform portability, applications are written in high-level languages such as Java. New systems and new algorithms can be put together quickly, but not optimal in terms of performance and power efficiency.

Our approach focusses on “drilling down” rather than scaling out. Storage, network, and computation should be better integrated to meet challenging system requirements. Judicious data processing in storage and directly off the network would significantly reduce transfer and activities on the host computer, leading to better performance and power efficiency. The architecture to realize this vision is shown below in Figure 1. A field-programmable gate-array (FPGA) is at the heart of the architecture. It performs three tasks: a) Control non-volatile memory (NVM) storage such as flash, b) Provide high-speed dedicated network with other storage modules, and c) Compute on data and interface to the processor(s) of the host computer. The whole unit operates as a plug-in accelerator to the computer. Many units can be connected to form a cluster, with as few or many computer “heads” as needed for the application. By tightly integrating storage, network, and computing together using an FPGA, up to a **10×** performance boost and **10×** power efficiency improvement could be achieved for application-specific designs.

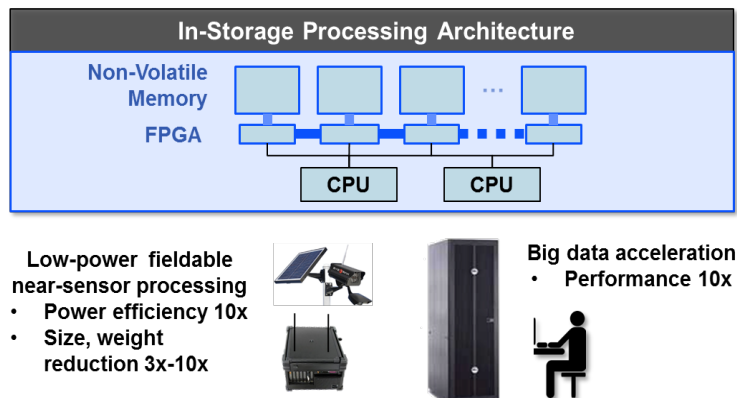


Figure 1. In-storage processing architecture.

This project team consists of research staff from MIT Lincoln Laboratory and MIT Computer Science and Artificial Intelligence Laboratory. Professor Arvind's group and an industry partner, Quanta Research, with additional support from Xilinx, Samsung, and Intel developed the architecture and prototype accelerator called Blue Database Management (BlueDBM) [1], consisting of custom-built flash modules with a local controller that interface with a Xilinx board hosting the FPGA. Software and firmware is custom designed for each application. Lincoln Laboratory provides support and guidance for system-level use cases and application concepts, which drive the development of high-level layer of software and firmware to expand the application space and make the accelerator useful for a wider community by lowering the barrier to entry. In addition, application insertion opportunities are being pursued within Lincoln Laboratory as well as external sponsors. Our vision is to bring new capabilities in big-data and internet-of-things applications with this new architecture.

Key research activities are listed below:

- A. Application-specific demos: These provide concrete results to illustrate potential benefits for classes of problems. In general, problems data sets larger than memory (1 TB) and random access pattern score very well. Software and firmware are often customized for each application at the C++/BlueSpec level. Example applications include:
 - a. Image similarity search
 - b. Document similarity search
 - c. Key-value store caching
 - d. Graph traversal
 - e. PageRank
 - f. Friends-of-friends
- B. Interface with open-source software: This allows BlueDBM technology to participate and contribute to the scale-out development approach used in the big-data algorithm prototyping community on LLGrid and beyond. A plug-and-play accelerator would provide out-of-the-box performance benefit, with further gains achieved through additional customization of the open-source software and BlueDBM wrapper. The interface to open-source software is most likely Java. Potential development and integration include:
 - a. D4M-on-BlueDBM graph traversal
 - b. Hadoop Distributed File System (HDFS) swap-in
 - c. Accumulo swap-in
- C. Size, weight, and power (SWaP) constrained fieldable system and new application concepts: Small fieldable units could provide "hotspot" storage, compute, and communication services for fielded sensors in remote or disaster areas. Potential tasks are:
 - a. Application studies
 - b. "Headless" BlueDBM pulls all control and processing from CPU into embedded processor in the FPGA

PROJECT TIMELINE

FY14 was focused on building the flash modules and developing firmware for the FPGAs to control the flash, provide high-speed networking, and interface to the computer. A few early benchmarks and system-level analysis were performed to demonstrate potential benefits. High-level database applications were surveyed and D4M was selected as a pilot package to be accelerated. [2]

In FY15, three mini applications were built (via BlueSpec) for demonstrations: a) Image similarity search, b) Key-value store cache for database queries, and c) Graph traversal on web-link dataset. The applications show a $7\times$ to $13\times$ performance improvement compared to standard server computer. Power efficiency is improved by $2\times$ to $5\times$. Two conference papers were presented in 2015 [3][4]. For open-source efforts, the software stacks of processing (D4M, Matlab) and database (Accumulo, HDFS) were analyzed for BlueDBM integration. Another open-source software package, Spark, was also selected as a candidate for acceleration. Spark has a large user community and rich algorithmic development environment spanning analytics, machine learning, etc.

In FY16, firmware and software infrastructure will be built up for BlueDBM to HDFS interface to support open-source software such as Accumulo and Spark. Two applications will be demonstrated to show performance gains and power savings. Application concepts exploration and follow-on sponsorship search activities will leverage FY15 developments.

This page intentionally left blank.

TABLE OF CONTENTS

	Page
Executive Summary	iii
List of Illustrations	ix
List of Tables	xi
1. INTRODUCTION	1
2. APPLICATION AND TECHNOLOGY TRENDS	3
3. TECHNICAL APPROACHES	5
3.1 BlueDBM Architecture	5
3.2 Open-System Architecture Considerations	7
3.3 Low-Power Application Concepts	9
4. FY15 ACCOMPLISHMENTS	11
4.1 Application-Specific Demos	11
4.2 Open-Source Software Interface Analysis	14
4.3 Low-Power Fieldable System and Application Concepts	15
5. SUMMARY	17
6. REFERENCES	19
APPENDIX A	21
APPENDIX B	35

This page intentionally left blank.

LIST OF ILLUSTRATIONS

Figure No.		Page
1	In-storage processing architecture.	iii
2	Uses of big-data analysis.	1
3	Data handling and processing pipeline.	3
4	“Hotspot” storage-compute-communicate services for remote sensors.	4
5	Traditional computer and power profile [10].	5
6	Previous approaches to boost performance.	6
7	BlueDBM approach.	6
8	BlueDBM software stack.	7
9	Application-level software interface to BlueDBM filesystem.	8
10	Optimization of application software for high performance.	8
11	Application software interface to BlueDBM accelerator	9
12	String search kernel.	11
13	Image similarity application.	12
14	Key-value caching.	13
15	One-server and four-server configurations.	14

This page intentionally left blank.

LIST OF TABLES

Figure No.		Page
1	Image Similarity Search	12
2	BlueDBM Advantage	13

This page intentionally left blank.

1. INTRODUCTION

The analysis of previously unimaginable amount of data can provide deep insight and expose certain information. Big data applications often involve the analysis of unstructured data, natural language text, videos and images, using advanced machine learning techniques. There are numerous examples of benefits from big data analysis. Google has predicted flu outbreaks a week earlier than the Center for Disease Control. Analyzing personal genomes can determine predisposition to diseases. Social network chatter analysis can identify political revolutions before newspapers. Scientific datasets can be mined to extract accurate models. Figure 2 shows that as data models become more sophisticated, we can progress from questions such as what happened, to why it happened, and what will happen, and how to make it happen.

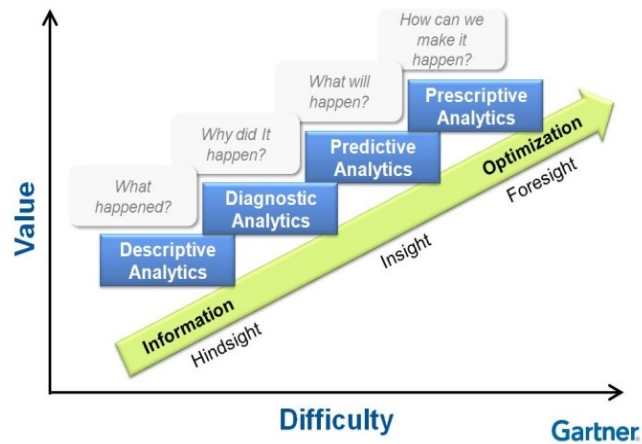


Figure 2. Uses of big-data analysis.

Data becomes so large and complex, growing 16x over the past 5 years, that traditional processing methods are no longer effective. The volume to be analyzed is too large to fit reasonably in DRAM computer memory. Furthermore, access is intensively random, making prefetching and caching ineffective. Data is often stored in the secondary storage of multiple machines on a cluster, thus, storage system and network performances become first-order concerns. Adding more computers brings only a diminishing return. Therefore, a new paradigm is needed.

In this project, we are exploring a novel near-data processing architecture through a collaborative effort with MIT CSAIL Professor Arvind. As depicted in Figure 1, the architecture is a plug-in accelerator offering three functions: a fast storage with embedded computation capability, a high-speed dedicated network with computation capability, and an efficient interface that interacts with the software on the

computer. All three functions are performed by a field-programmable gate-array (FPGA), which can provide superior performance and power efficiency compared to a general purpose computer processor. An improvement of $3\times$ – $10\times$ in performance and power efficiency could be achieved. The scope of applications include back-office analytics acceleration, as well as field providing field storage-compute-communication services to sensors in remote or disaster relief areas.

The rest of the report will discuss application and technology trends, followed by the technical approach, FY15 accomplishments, and future work. Two conference papers [3][4] were published on the accelerator architecture and mini-app performances in FY15. This report will address unpublished aspects of the project.

2. APPLICATION AND TECHNOLOGY TRENDS

Big data applications involve structure as well as unstructured data. New types of analysis include text analytics for topic modeling, image similarity matching, face recognition, audio and video labeling, etc. The analytics processing typically follows a data handling and processing pipeline involving parsing, database ingesting, scanning and query and enriching a subset of interest, and finally performing analysis. Figure 3 illustrates a processing pipeline example where D4M is used as an analysis environment.

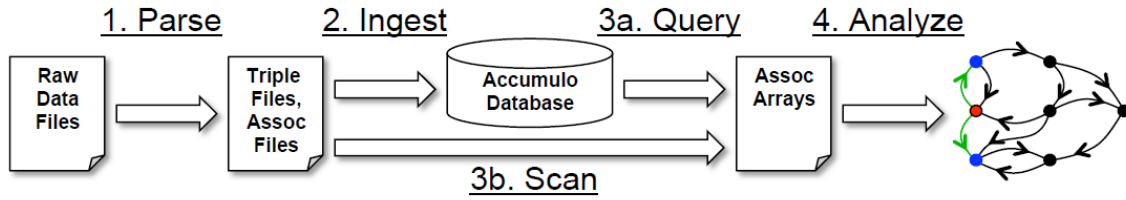


Figure 3. Data handling and processing pipeline.

D4M, or Dynamic Distributed Dimensional Data Model, is a Lincoln Laboratory–developed technology that offers rapid prototyping capabilities to scientists interested in working with large datasets [5]. D4M provides support for database and computation systems that deal with big data through the mathematical foundations of Associative Arrays, schema designed for NoSQL key-value store databases, and the ability to interact with diverse database technologies such as SQL, Accumulo [6] and SciDB. D4M is implemented with a MATLAB API. D4M has been adopted by a large number of users within and outside Lincoln Laboratory, largely due to its low barrier to entry, easy to adopt API, and applicability to a large variety of unstructured data.

Spark is another widely used open-source software package that provides an environment for analytics algorithm development. Spark [7] was developed at U.C. Berkeley before released into the Apache open-source domain. Its inventor has recently joined MIT faculty. Spark has a large industry user base, and offers a unified environment for different types of analytics processing, from database query to graph processing and machine learning.

Interestingly, both D4M/Accumulo and Spark use Hadoop Distributed File System (HDFS) as the database file system. HDFS [8] was developed by Google and is an integral part of the Map-Reduce distributed computing scheme that match up processing tasks to local data at each node – sending processing tasks to data rather than traditionally moving data around to processors. Another interesting point is all three packages, D4M, Accumulo, Spark, HDFS, are implemented with Java or a variant. Such languages can be a factor of 1/3 as slow as C/C++, but it allows for the application code to run on any platform without recompilation.

Clearly, the current state-of-the art approach for big data is scaling out to many computers to meet requirements on processing, storage capacity, and access bandwidth. With the emphasis on scalability, companies use commodity-class computers and rely on quantity and technology upgrades to meet performance. Systems and new algorithms can be put together quickly, but are not optimal in performance and power efficiency.

There are emerging applications that require all available processing available to meet desired user latency. These applications include natural language processing, voice and image query, and typical Apple or Google cloud processing. Special accelerators, such as GPUs and FPGAs, are being explored to help increase performance. It is estimated that a future system will have a mixture of processing technologies to meet performance and power budget [9].

Our accelerator will need to fit well in this ecosystem to make a substantial impact. On one hand, it needs to be work well with the open-source “scale-out” approach. It would need to conform to some open-source interfaces and protocols. On the other hand, it also needs to be customizable to take on application-specific acceleration, for example, in-storage processing or in an extreme case, low power embedded processing of sensor data in remote area.

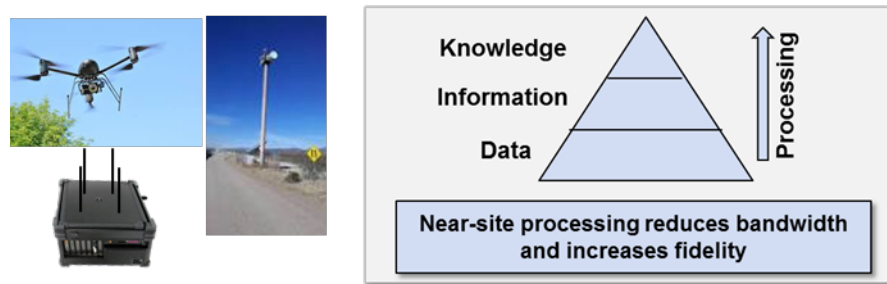


Figure 4. “Hotspot” storage-compute-communicate services for remote sensors.

3. TECHNICAL APPROACHES

In this section, we will describe our detailed approaches for the BlueDBM architecture and the open-system software architecture used for algorithm testbed and application development.

3.1 BLUEDBM ARCHITECTURE

Traditional computer systems pull data from disks into memory for subsequent processing. As datasets exceed available memory, portions must get swapped out to disks during the processing, resulting in large amount of disk transfer activity. Disk access is slow and also causes high workload for the computer CPU.

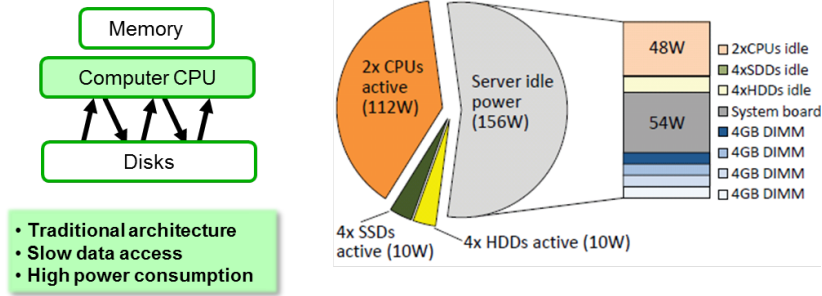


Figure 5. Traditional computer and power profile [10].

Recently, driven by big-data processing demand and falling price of DRAM, new systems incorporate more memory to reduce disk access [11]. To boost performance even further, Microsoft has started using FPGAs in the processing of Bing queries [12]. These approaches achieve significantly improved performance but with additional cost and power consumption.

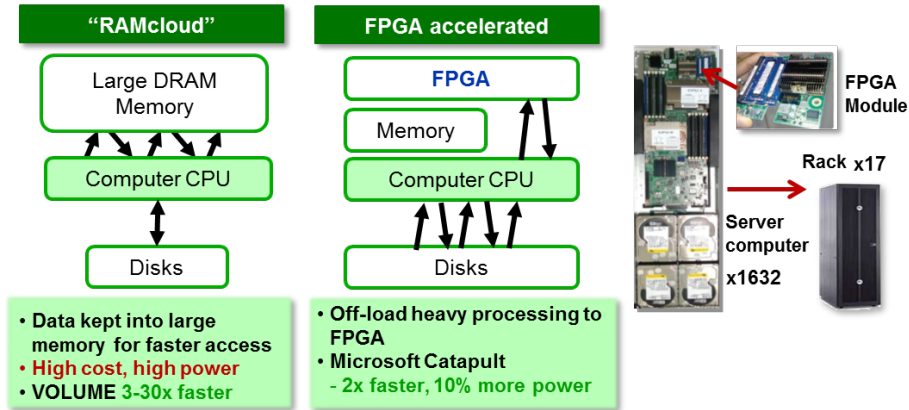


Figure 6. Previous approaches to boost performance.

Our BlueDBM architecture, depicted below, incorporates the FPGA as part of the storage system to perform computations, reducing the amount of data handling and processing by the computer processor. This leads to benefits both in performance and power consumption. Detail specifications are available in [2][3].

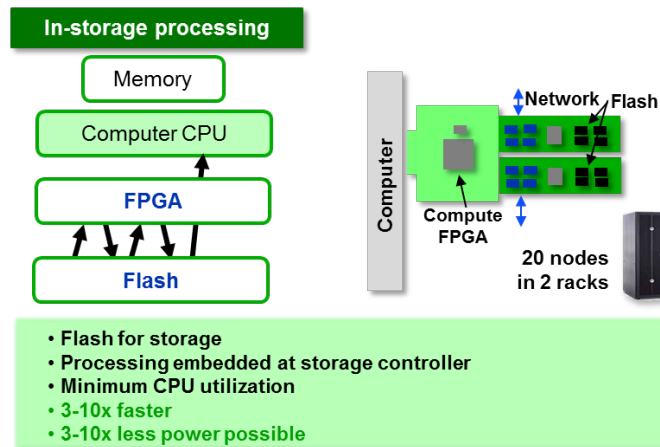


Figure 7. BlueDBM approach.

The accelerator is interfaced to the host computer through the BlueDBM software stack [3]. The stack allows for user-level access to files. For analytics acceleration, the user-level software would first obtain physical locations of the data, then dispatch queries and analysis tasks to the FPGA through a

special conduit. The results can be kept in the file system or sent back to the computer. User-level in this context refers to the BlueSpec environment, which frees the user from low-level VHDL coding, but still requires synthesis and C/C++ level programming. This approach is best used for developing customized functions.

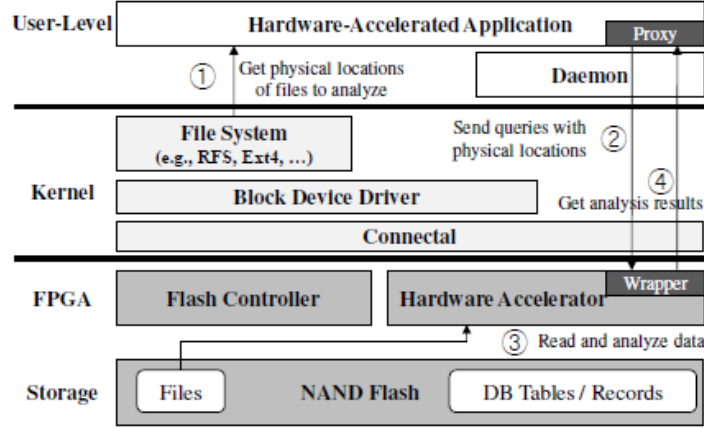


Figure 8. BlueDBM software stack.

3.2 OPEN-SYSTEM ARCHITECTURE CONSIDERATIONS

As discussed earlier, many open-source database and analytics software packages are architected to be cross-platform portable. Accumulo/D4M is used for sparse database and analysis, whereas Spark operates in different domains such as structured and streaming data. Accumulo and Spark both interface with HDFS, which provides file management services at the Java level. In order for BlueDBM to support such packages, it will need an HDFS wrapper, so that it can be accessed seamlessly. This is depicted in Figure 9 below. Blue represents existing code with no modification.

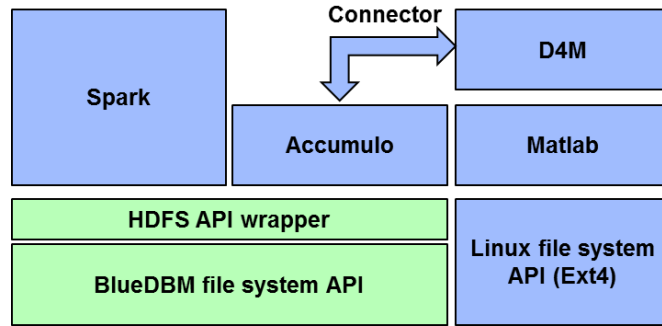


Figure 9. Application-level software interface to BlueDBM filesystem.

To fully benefit from the performance offered by BlueDBM, the packages may need some modification to improve data transport mechanisms. For example, if file access is performed through an inefficient software stack, then the speed up underneath is diminished by the slow API. Figure 10 depicts a scenario with higher performance. The required modification is expected to be minor.

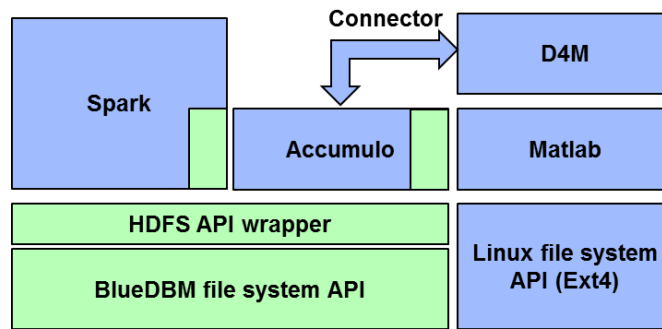


Figure 10. Optimization of application software for high performance.

On the analytic side, acceleration needs to go through the steps outlined in Figure 8 for the BlueDBM software stack. Since the process involves calls through user space and kernel space, it is best implemented as a manager unit and presented up to application software via an accelerator API. In the case of Accumulo/D4M, this API needs to interface with Java code. The accelerator API could take on a different flavor pending on the application software; it could be an interface for D4M when it is used alone without Accumulo. Within D4M or Accumulo, additional development will be needed to incorporate the acceleration capability. The architecture for integration with BlueDBM accelerator is depicted below.

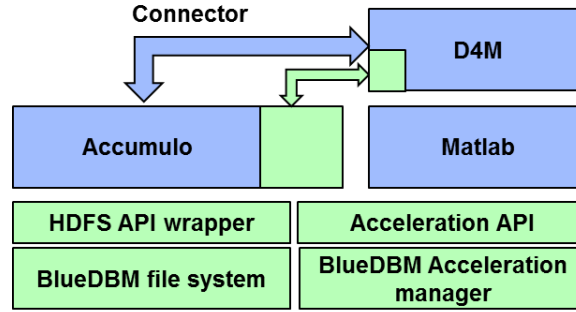


Figure 11. Application software interface to BlueDBM accelerator

In terms of implementation, the platform would start with a single server with an option to scale up and be part of Lincoln Laboratory’s LLGrid system.

3.3 LOW-POWER APPLICATION CONCEPTS

In addition to improving performance, off-loading activities from the processor to the FPGA also results in better power efficiency. By tightly integrating storage, network, and computing together using an FPGA, up to a **10×** performance boost and **10×** power efficiency improvement could be achieved for application-specific designs. Notably, the reduction in power would open up new capabilities such as off-grid portable field operation for storage-compute-communication services.

In general, low-power operation often has peaks and troughs, where system-level optimizations such as clock gating, reduced duty cycle operation, etc. could bring out further power savings. For “cloud-like” field services, the processing would need to be performed by the FPGA [9], the processor footprint in the system could be further reduced into the embedded realm, or even pulled completely inside the FPGA if appropriate.

This page intentionally left blank.

4. FY15 ACCOMPLISHMENTS

FY15 was focused on the following activities: application-specific demonstrations, open-source software interface analysis, and low-power fieldable system and application concepts.

4.1 APPLICATION-SPECIFIC DEMOS

The application-specific demos provide illustrate potential benefits for certain classes of problems. In general, challenging problems with large data sets larger than memory (1 TB) and random access pattern score really well. Software and firmware are often customized for each application at the C++/BlueSpec level.

Listed below are the kernel and applications. The results are in the attached publications [3][4], so only selected highlights are discussed in the report.

1. String search kernel
2. Image similarity search application
3. Document similarity search application
4. Key-value store caching application
5. Graph traversal application
6. PageRank application (FY16)
7. Friends-of-friends application (FY16)

String search is a common kernel in big-data applications, spanning from DNA sequencing to text analytics. A string search engine was built on BlueDBM using Morris-Pratt algorithm. Figure 12 shows BlueDBM in-storage processing yields **7.5x** boost relative to CPU+disk, where the bottleneck is storage. When faster storage – flash is used, performance increases with high CPU activity. Off-loading to FPGA reduces CPU activity by **20x**.

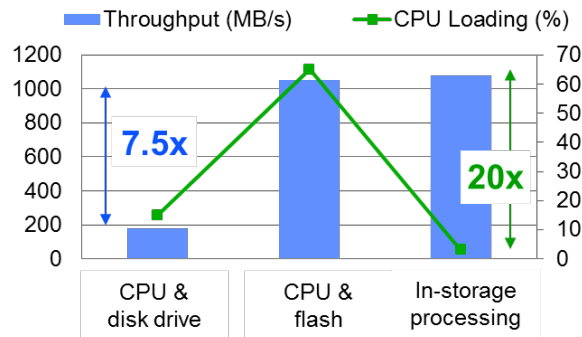


Figure 12. String search kernel.

Illustrated in Figure 13 is an image similarity search application. The search subjects an image to a large dataset of 80 million images (~300 GB) to return a set of similar images, their labels, and match metrics. To reduce workload, the image set is organized (indexed) by histogram similarity. Thus, the search is narrowed down to subsets prior to a thorough match is performed. The images in each indexed subset are scattered in physical storage.

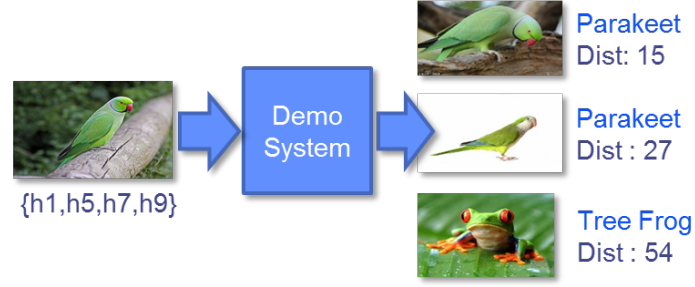


Figure 13. Image similarity application.

Performance results are shown in Tables 1 and 2 below. Relative to baseline configuration using CPU+hard drives, copying the dataset (in this case 30 GB) into DRAM and running the search on the CPU shows 9.2x performance boost. Note that this does not scale up for large datasets. BlueDBM with FPGA+flash yields a speed up **12x** for dataset up to 1 TB per server and larger using more computers.

TABLE 1
Image Similarity Search

Configuration	Performance	Power
CPU + hard drives	1x reference	140W = 90W idle + 50W activity
CPU + DRAM	9.2x speedup	240W = 90W idle + 150W activity
BlueDBM FPGA + flash	12x speed up	120W = 90W idle + 30W activity

Power-wise, BlueDBM offers peak performance at 30W activity power (or 120W if including computer system idle power, which could be reduced significantly for application specific scenario). BlueDBM activity power is only **1/5** compared to CPU+DRAM configuration.

TABLE 2
BlueDBM Advantage

Configuration	Performance	Power
BlueDBM vs CPU+hard drives	12x speedup	50/30 = 1.6x less activity power
BlueDBM vs CPU+DRAM	1.3x speedup	150/30 = 5x less activity power

Key-value query is another form of search geared toward unstructured data. A page update on Amazon or Facebook generates tens to hundreds of queries. The keys and values are cached in the memory of web servers to provide faster response compared to reading from file systems. The cache is distributed over many web servers, so that a miss from one could be serviced by another over a computer network. Computer memory, however, is expensive compared to disks and flash. Furthermore, computer network has long latency relative to memory access, so the benefit of distributed memory is limited.

BlueDBM uses flash rather than DRAM memory for caching, its FPGA-based network rather than a generic computer network, and FPGA logic to perform caching. The integration of storage + network + compute functions into one FPGA also makes the flow well streamlined. Figure 14 depicts the use of BlueDBM in key-value caching. Compared to the CPU/DRAM configuration, the BlueDBM FPGA/flash/ network solution is **13x** faster, not to mention less expensive and lower power.

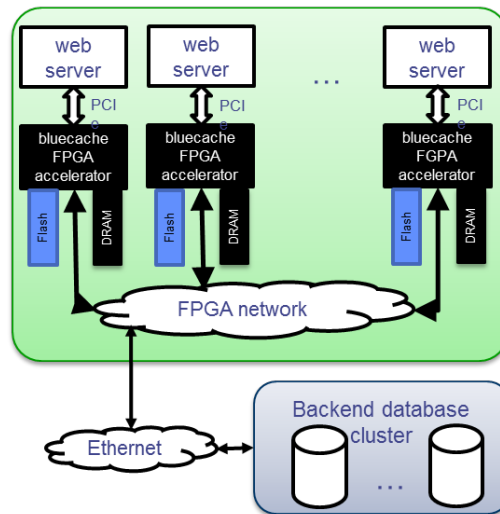


Figure 14. Key-value caching.

4.2 OPEN-SOURCE SOFTWARE INTERFACE ANALYSIS

Interfacing with open-source software would allow BlueDBM technology to participate and contribute to the scale-out development approach used in big-data algorithm prototyping community. After analyzing Accumulo and Spark open-source software packages, we arrived at a vision for BlueDBM usage in this environment: A plug-and-play accelerator would provide out-of-the-box performance benefit, with further gains achieved through additional customization of the open-source software and BlueDBM wrapper. The software architecture approach is described in Section 3.2, Open-System Architecture Considerations.

We've identified three potential applications using BlueDBM with open-source software:

1. D4M/Accumulo graph traversal on BlueDBM
 - Runs existing algorithms, with minor modification
 - Enables out-of-core data size for D4M
 - Accelerates sparse matrix multiplication in FPGA
2. Hadoop Distributed File System (HDFS) swap-in
 - BlueDBM with HDFS look-and-feel
 - Supports many software packages on LLGrid, including Accumulo and Spark
 - Exceeds maximum throughput of HDFS client software
3. Accumulo swap-in
 - BlueDBM key-value store with Accumulo look-and-feel
 - Supports many software packages, including D4M
 - Accelerates database operation
 - Extreme performance with swap-in HDFS

The platform starts with a single server with options to scale up to four servers and be part of Lincoln Laboratory's LLGrid system.



Figure 15. One-server and four-server configurations.

4.3 LOW-POWER FIELDABLE SYSTEM AND APPLICATION CONCEPTS

We've explored a few application concepts for low-power fieldable system. One notable concept is to provide "hotspot" storage, compute, and communication services for fielded sensors in remote or disaster areas. Enabled by compute capability, the system could also provide "cloud-like" services for local users.

Through discussions with colleagues, we have identified open-domain large datasets for big-data applications that can be shared for collaboration. A set of examples is listed below:

- <https://www.data.gov>
- <http://www.ncdc.noaa.gov/data-access/quick-links#loc-clim>
- http://oad.simmons.edu/oadwiki/Data_repositories
- <https://library.uoregon.edu/datamanagement/repositories.html>
- <http://www.nature.com/sdata/data-policies/repositories>
- <http://aws.amazon.com/datasets>
 - **Size:** 200 GB, **Source:** The US Census Bureau
 - **Size:** 210 GB, **Source:** Fed. Energy Reg. Commission (FERC) investigate Enron
 - **Size:** 500 GB, **Source:** <http://labrosa.ee.columbia.edu/millionsong/>
 - **Size:** 800 GB, **Source:** Data Wrangling – wiki traffic
 - **Size:** 2.2 TB, **Source:** Google Books
 - **Size:** 5 TB, **Source:** Human Microbiome Project
 - **Size:** 5 TB, **Submitted By:** modENCODE DCC (help@modencode.org)
 - **Size:** 6 TB, **Source:** CCAFS Climate Portal www.ccafs-climate.org
 - **Size:** 200 TB, **Source:** National Center for Biotechnology Information (NCBI)
 - **Size:** 541 TB, **Source:** Common Crawl Foundation - <http://commoncrawl.org>

This page intentionally left blank.

5. SUMMARY

The development of a low-power embedded analytics accelerator is a three-year university collaboration between Lincoln Laboratory and Professor Arvind's group at the MIT Computer Science and AI Laboratory (CSAIL). The goal of the project is to design and prototype a novel architecture that has wide potential applicability to the new types of applications ranging from back-office big-data analytics to fieldable hot-spot systems providing storage-processing-communication services for off-grid sensors. Our approach simultaneously optimizes storage, network, and computation performance within the system. Up to **10×** performance boost and **10×** power efficiency improvement could be achieved for application-specific designs. Our vision is to bring new capabilities in big-data and internet-of-things applications with this new architecture.

Key research activities are listed below:

- A. Application-specific demos
- B. Interface with open-source software
- C. Size, weight, and power (SWaP) constrained fieldable system and application concepts

FY14 was focused on building the BlueDBM prototype. In FY15, three mini applications were built (via BlueSpec) for demonstrations. The applications show **7×** - **13×** in performance compared to a standard server computer, and power efficiency is improved by **2×** to **5×**. With respect to the open-source actions, the software stacks for D4M, Accumulo, HDFS, Spark were investigated for BlueDBM integration. In FY16, firmware and software infrastructure will be built up for BlueDBM to interface with open-source software. Application concepts exploration and follow-on sponsorship search activities will leverage FY15 developments.

This page intentionally left blank.

6. REFERENCES

- [1] Arvind, “BlueDBM: A Multi-access, Distributed Flash Store for Big Data Analytics,” Keynote presentation, IEEE High Performance Embedded Computing, September 2014
- [2] H. Nguyen, V. Gadepally, J. Muldavin, and Arvind, “Low-Power Embedded Analytics: FY14 Line-Supported Information, Computation, and Exploitation Program”, Project Report LSP-124
- [3] S. Jun, M. Liu, S. Lee, J. Hicks, J. Ankcorn, M. King, S. Xu, Arvind, “BlueDBM: An Appliance for Big Data Analytics,” Intl Symposium on Computer Architecture (ISCA), Portland, OR, June 2015
- [4] S. Jun, M. Liu, S. Xu, Arvind, “A Transport-Layer Network for Distributed FPGA Platforms, International Conference on Field-programmable Logic and Applications (FPL),” London, UK, September 2015
- [5] J. Kepner et al., “Dynamic Distributed Dimensional Data Model (D4M) Database and Computation System,” ICASSP, March 2012
- [6] M. Wall, A. Cordova and B. Rinaldi, Accumulo Application Development, Table Design, and Best Practices, O’Reilly, Sebastapol, California, US, 2013
- [7] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M. J. Franklin, S. Shenker, I. Stoica, “Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing,” NSDI’12
- [8] A. Bialecki et al., “Hadoop: A Framework for Running Applications on Large Clusters Built of Commodity Hardware,” 2005. Wiki at <http://lucene.apache.org/hadoop>.
- [9] J. Hauswald, M. A. Laurenzano, Y. Zhang, C. Li, A. Rovinski, A. Khurana, R. Dreslinski, T. Mudge, V. Petrucci, L. Tang, and J. Mars, “Sirius: An Open End-to-End Voice and Vision Personal Assistant and Its Implications for Future Warehouse Scale Computers.” In Proceedings of the Twentieth International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), ASPLOS ’15, New York, NY, USA, 2015. ACM
- [10] Tsirogiannis, Harizopoulos, Shah, “Analyzing the Energy Efficiency of a Database Server,” SIGMOD, 2010
- [11] Ma, Hong, Gu, “VOLUME: Enable arge-Scale In-Memory Computation on Commodity Cluster,” CloudCom’13
- [12] Putnam et al., “A Reconfigurable Fabric for Accelerating Large-Scale Datacenter Services”, ISCA, 2014

This page intentionally left blank.

APPENDIX A

BlueDBM: An Appliance for Big Data Analytics

Sang-Woo Jun[†] Ming Liu[†] Sungjin Lee[†] Jamey Hicks^{*}
John Ankcorn^{*} Myron King^{*} Shuotao Xu[†] Arvind[†]

Department of Electrical Engineering and Computer Science
Massachusetts Institute of Technology[†]
Quanta Research Cambridge^{*}

{wjun, ml, chamdoo, shuotao, arvind}@csail.mit.edu[†]
{jamey.hicks, john.ankcorn, myron.king}@qrclab.com^{*}

Abstract

Complex data queries, because of their need for random accesses, have proven to be slow unless all the data can be accommodated in DRAM. There are many domains, such as genomics, geological data and daily twitter feeds where the datasets of interest are 5TB to 20 TB. For such a dataset, one would need a cluster with 100 servers, each with 128GB to 256GBs of DRAM, to accommodate all the data in DRAM. On the other hand, such datasets could be stored easily in the flash memory of a rack-sized cluster. Flash storage has much better random access performance than hard disks, which makes it desirable for analytics workloads. In this paper we present BlueDBM, a new system architecture which has flash-based storage with in-store processing capability and a low-latency high-throughput inter-controller network. We show that BlueDBM outperforms a flash-based system without these features by a factor of 10 for some important applications. While the performance of a ram-cloud system falls sharply even if only 5%~10% of the references are to the secondary storage, this sharp performance degradation is not an issue in BlueDBM. BlueDBM presents an attractive point in the cost-performance trade-off for Big Data analytics.

1. Introduction

By many accounts, complex analysis of Big Data is going to be the biggest economic driver for the IT industry. For example, Google has predicted flu outbreaks by analyzing social network information a week faster than CDC [13]; Analysis of twitter data can reveal social upheavals faster than journalists; Amazon is planning to use customer data for anticipatory shipping of products [43]; Real-time analysis of personal genome may significantly aid in diagnostics. Big Data analytics are

potentially going to have revolutionary impact on the way scientific discoveries are made.

Big Data by definition doesn't fit in personal computers or DRAM of even moderate size clusters. Since the data may be stored on hard disks, latency and throughput of storage access is of primary concern. Historically, this has been mitigated by organizing the processing of data in a highly sequential manner. However, complex queries cannot always be organized for sequential data accesses, and thus high performance implementations of such queries pose a great challenge. One approach to solving this problem is *ram cloud* [34], where the cluster has enough collective DRAM to accommodate the entire dataset in DRAM. In this paper, we explore a much cheaper alternative where Big Data analytics can be done with reasonable efficiency in a single rack with distributed flash storage, which has much better random accesses performance than hard disks. We call our system BlueDBM and it provides the following capabilities:

1. A 20-node system with large enough flash storage to host Big Data workloads up to 20 TBs;
2. Near-uniform latency access into a network of storage devices that form a global address space;
3. Capacity to implement user-defined in-store processing engines;
4. Flash card design which exposes an interface to make application-specific optimizations in flash accesses.

Our preliminary experimental results show that for some applications, BlueDBM performance is an order of magnitude better than a conventional cluster where SSDs are used only as a disk replacement. BlueDBM unambiguously establishes an architecture whose price-performance-power characteristics provide an attractive alternative for doing similar scale applications in a ram cloud.

As we will discuss in the related work section, almost every element of our system is present in some commercial system. Yet our system architecture as a whole is unique. The main contributions of this work are: (1) Design and implementation of a scalable flash-based system with a global address space, in-store computing capability and a flexible inter-controller network. (2) A hardware-software codesign environment for

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage, and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ISCA'15, June 13-17, 2015, Portland, OR USA

Rights management text and bibliographic strip from ACM placed here.

incorporating user-defined in-store processing engines. (3) Performance measurements that show the advantage of such an architecture over using flash as a drop-in replacement for disks. (4) Demonstration of a complex data analytics appliance which is much cheaper and consumes an order of magnitude less power than the cloud-based alternative.

The rest of the paper is organized as follows: In Section 2 we explore some existing research related to our system. In Section 3 we describe the architecture of our rack-level system, and in Section 4 we describe the software interface that can be used to access flash and the accelerators. In Section 5 we describe a hardware implementation of BlueDBM, and show our results from the implementation in Section 6. In Section 7 we describe and evaluate some example accelerators we have built for the BlueDBM system. Section 8 summarizes our paper.

2. Related Work

In Big Data scale workloads, building a cluster with enough DRAM capacity to accommodate the entire dataset can be very desirable but expensive. An example of such a system is RAMCloud, which is a DRAM-based storage for large-scale datacenter applications [34, 39]. RAMCloud provides more than 64TBs of DRAM storage distributed across over 1000 servers networked over high-speed interconnect. Although RAMCloud provides 100 to 1000 times better performance than disk-based systems of similar scale, its high energy consumption and high price per GB limits its widespread use except for extremely performance and latency-sensitive workloads.

NAND-Flash-based SSD devices are gaining traction as a faster alternative to disks, and close the performance gap between DRAM and persistent storage. SSDs are an order of magnitude cheaper price compared to DRAM, and an order of magnitude faster performance compared to disk. Many existing database and analytics software has shown improved performance with SSDs [8, 21, 27]. Several SSD-optimized analytics softwares, such as the SanDisk Zetascale [40] have demonstrated promising performance while using SSD as the primary data storage. Many commercial SSD devices have adopted high-performance PCIe interface in order to overcome the slower SATA bus interface designed for disk [11, 30, 16]. Attempts to use flash as a persistent DRAM alternative by plugging it into a RAM slot are also being explored [45].

SSD storage devices have been largely developed to be a faster drop-in replacement for disk drives. This backwards compatibility has helped their widespread adoption. However, additional software and hardware is required to hide the difference in device characteristics [1]. Due to the high performance of SSDs, even inefficiencies in the storage management software becomes significant, and optimizing such software has been under active investigation. Moneta [4] modifies the operating system's storage management components to reduce software overhead when accessing NVM storage devices. Wil-

low [41] provides an easy way to augment SSD controllers with additional interface semantics that make better use of SSD characteristics, in addition to a backwards compatible storage interface. Attempts to remove the translation layers and let the database make high-level decisions [14] have shown to be beneficial.

Due to their high performance, SSDs also affect the network requirements. The latency to access disk over Ethernet was dominated by the disk seek latency. However, in a SSD-based cluster the storage access latency could even be lower than network access. These concerns are being addressed by faster network fabrics such as 10GbE and Infiniband [2], and by low-overhead software protocols such as RDMA [29, 17, 38, 46, 29, 37] or user-level TCP stacks that bypass the operating system [19, 15]. QuickSAN [5] is an attempt to remove a layer of software overhead by augmenting the storage device with a low-latency NIC, so that remote storage access does not need to go through a separate network software stack.

Another important attempt to accelerate SSD storage performance is in-store processing, where some data analytics is offloaded to embedded processors inside SSDs. These processors have extremely low-latency access to storage, and helped overcome the limitations of the storage interface bus. The idea of in-store processing itself is not new. Intelligent disks (IDISK) connected to each other using serial networks have been proposed in 1998 [23], and adding processor to disk heads to do simple filters have been suggested as early as in the 1970s [28, 35, 3]. However, performance improvements of such special purpose hardware did not justify their cost at the time.

In-store processing is seeing new light with the advancement of fast flash technology. Devices such as Smart SSDs [9, 22, 41] and Programmable SSDs [6] have shown promising results, but gains are often limited by the performance of the embedded processors in such power constrained devices. Embedding reconfigurable hardware in storage devices is being investigated as well. For example, Ibex [48] is a MySQL accelerator platform where a SATA SSD is coupled with an FPGA. Relational operators such as selection and group-by are performed on the FPGA whenever possible, otherwise they are forwarded to software. Companies such as IBM/Netezza [42] offload operations such as filtering to a reconfigurable fabric near storage. On the other end of the spectrum, systems such as XSD [6] embeds a GPU into a SSD controller, and demonstrates high performance accelerating MapReduce.

Building specialized hardware for databases have been extensively studied and productized. Companies such as Oracle [33] have used FPGAs to offload database queries. FPGAs have been used to accelerate operations such as hash index lookups [25]. Domain-specific processors for database queries are being developed [44, 47], including Q100 [49] and LINQits [7]. Q100 is a data-flow style processor with an instruction set architecture that supports SQL queries. LINQits

mapped a query language called LINQ to a set of accelerated hardware templates on a heterogeneous SoC (FPGA + ARM). Both designs exhibited order of magnitude performance gains at lower power, affirming that specialized hardware for data processing is very advantageous. However, unlike BlueDBM, these architectures accelerate computation on data that is in DRAM. Accelerators have also been placed in-path between network and processor to perform operations at wire speed [32], or to collect information such as histogram tables without overhead [18].

Incorporating reconfigurable hardware accelerators into large datacenters is also being investigated actively. Microsoft recently has built and demonstrated the power/performance benefits of an FPGA-based system called Catapult [36]. Catapult uses a large number of homogeneous servers each augmented with an FPGA. The FPGAs form a network among themselves via high-speed serial links so that large jobs can be mapped to groups of FPGAs. Catapult was demonstrated to deliver much faster performance while consuming less power, compared to a normal ram cloud cluster. BlueDBM has similar goals in terms of reconfigurable hardware acceleration, but it uses flash devices to accelerate lower cost systems that do not have enough collective DRAM to host the entire dataset.

This system improves upon our previous BlueDBM prototype [20], which was a 4-node system with less than 100GB of slow flash. It was difficult to extrapolate the performance of real applications from the results obtained from our previous prototype, because of both its size and different relative performance of various system components. The current generation of BlueDBM has been built with the explicit goal of running real applications, and will be freely available to the community for developing Big Data applications.

3. System Architecture

The BlueDBM architecture is a homogeneous cluster of host servers coupled with a BlueDBM storage device (See Figure 1). Each BlueDBM storage device is plugged into the host server via a PCIe link, and it consists of flash storage, an in-store processing engine, multiple high-speed network interfaces and on-board DRAM. The host servers are networked together using Ethernet or other general-purpose networking fabric. The host server can access the BlueDBM storage device via a host interface implemented over PCIe. It can either directly communicate with the flash interface, to treat it as a raw storage device, or with the in-store processor to perform computation on the data.

The in-store processing engine has access to four major services: The flash interface, network interface, host interface and the on-storage DRAM buffer. Figure 1 and Figure 2 shows the four services available to the in-store processor. In the following sections we describe the flash interface, network interface and host interface in order. We omit the DRAM buffer because there is nothing special about its design.

3.1. Flash Interface

Flash devices or SSDs achieve high bandwidth by grouping multiple flash chips into several channels, all of which can operate in parallel. Because NAND flash has limited program/erase cycles and frequent errors, complex flash management algorithms are required to guarantee reliability. These include wear leveling, garbage collection, bit error correction and bad block management. These functions are typically handled by multiple ARM-based cores in the SSD controller. The host side interface of an SSD is typically SATA or PCIe, using AHCI or NVMe protocols to communicate with host. SSDs are viewed as a typical block device to the host operating system, and its internal architecture and management algorithms are completely hidden.

However, this additional layer of management duplicates some file system functions and adds significant latency [26]. Furthermore, in a distributed storage environment, such as BlueDBM, independent flash devices do not have a holistic view of the system and thus cannot efficiently manage flash. Finally, in-store processors that we have introduced in BlueDBM would also incur performance penalties if passing through this extra layer. Thus in BlueDBM, we chose to shift flash management away from the device and into file system/block device driver (discussed in Section 4).

3.1.1. Interface for High Performance Flash Access Our flash controller exposes a low-level, thin, fast and bit-error corrected hardware interface to raw NAND flash chips, buses, blocks and pages. This has the benefit of (i) cutting down on access latency from the network and in-store processors; (ii) exposing all degrees of parallelism of the device and (iii) allowing higher level system stacks (file system, database storage engine) to more intelligently manage data.

To access the flash, the user first issues a flash command with the operation, the address and a tag to identify the request. For writes, the user then awaits for a write data request from the controller scheduler, which tells the user that the flash controller is ready to receive the data for that write. The user will send the write data corresponding to that request in 128-bit bursts. The controller returns an acknowledgement once write is finished. For read operations, the data is returned in 128-bit bursts along with the request tag. For maximum performance, the controller may send these data bursts *out of order* with respect to the issued request and *interleaved* with other read requests. Thus completion buffers may be required on the user side to maintain FIFO characteristics. Furthermore, we note that to saturate the bandwidth of the flash device, multiple commands must be in-flight at the same time, since flash operations can have latencies of 50 μ s or more.

3.1.2. Multiple Access Agents Multiple hardware endpoints in BlueDBM may need shared access to this flash controller interface. For example, a particular controller may be accessed by local in-store processors, local host software over PCIe

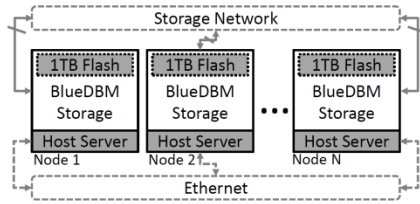


Figure 1: BlueDBM Overall Architecture

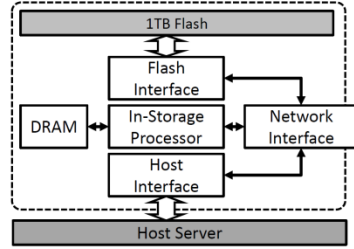


Figure 2: BlueDBM Node Architecture

DMA, or remote in-store processors over the network. Thus we implemented a Flash Interface Splitter with tag renaming to manage multiple users (Figure 3). In addition, to ease development of hardware in-store processors, we also provide an optional Flash Server module as part of BlueDBM. This server converts the out-of-order and interleaved flash interface into multiple simple in-order request/response interfaces using page buffers. It also contains an Address Translation Unit that maps file handles to incoming streams of physical addresses from the host. The in-store processor simply makes a request with the file handle, offset and length, and the Flash Server will perform the flash operation at the corresponding physical location. The software support for this function is discussed in Section 4). The Flash Server's width, command queue depth and number of interfaces is adjustable based on the application.

3.2. Integrated Storage Network

BlueDBM provides a low-latency high-bandwidth network infrastructure across all BlueDBM storage devices in the cluster, using a simple design with low buffer requirements. BlueDBM storage devices form a separate network among themselves

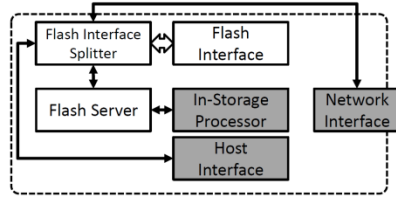


Figure 3: Flash Interface

via high-performance serial links. The BlueDBM network is a packet-switched mesh network, in which each storage device has multiple network ports and is capable of routing packets across the network without requiring a separate switch or router. In addition to routing, the storage network supports functionality such as flow control and virtual channels while maintaining high performance and extremely low latency. For data traffic between the storage devices, the integrated network ports removes the overhead of going to the host software to access a separate network interface.

Figure 4 shows the network architecture. Switching is done at two levels, the internal switch and the external switch. The internal switch routes packets between local components. The external switch accesses multiple physical network ports, and is responsible for forwarding data from one port to another in order to relay a packet to its next hop. It is also responsible for relaying inbound packets to the internal switch, and relaying outbound packets from the internal switch to a correct physical port.

Due to the multiple ports on the storage nodes, the BlueDBM network is very flexible and can be configured to implement various topologies, as long as there is sufficient number of ports on each node. Figure 5 shows some example topologies. To implement a different topology the physical cables between each node has to be re-wired, but the routing across a topology can be configured dynamically by the software.

3.2.1. Logical Endpoint The BlueDBM network infrastructure exposes virtual channel semantics to the users of the network by providing it with multiple *logical endpoints*. The number of endpoints are determined at design time by setting a parameter, and all endpoints share the physical network. Each endpoint is parameterized with a unique index that does not need to be contiguous. Each endpoint exposes two interfaces, *send* and *receive*. An in-store processor can send data to a

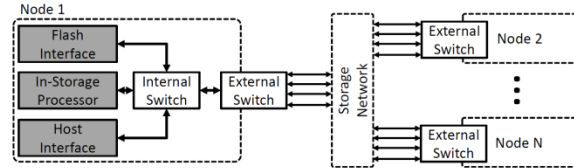


Figure 4: Network Architecture

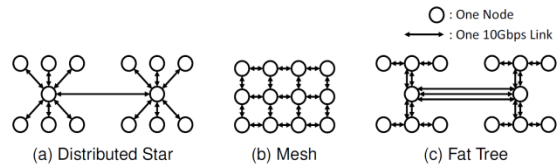


Figure 5: Any Network Topology Is Possible As Long As It Requires Less Than 8 Network Ports Per Node

remote node by calling `send` with a pair of data and destination node index, or receive data from remote nodes by calling `receive`, which returns a pair of data and source node index. These interfaces provide back pressure, so that each endpoint can be treated like a FIFO interface across the whole cluster. Such intuitive characteristics of the network ease development of in-store processors.

3.2.2. Link Layer The link layer manages physical connections between network ports in the storage nodes. The most important aspect of the link layer is the simple token-based flow control implementation. This provides back pressure across the link and ensures that packets will not drop if the data rate is higher than what the network can manage, or if the data cannot be received by the destination node which is running slowly.

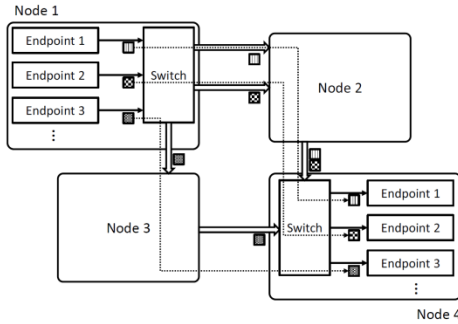


Figure 6: Packets From The Same Endpoint to a Destination Maintain FIFO Order

3.2.3. Routing Layer In order to make maximum use of the bandwidth of the network infrastructure while keeping resource usage to a minimal, the BlueDBM network implements deterministic routing for each logical endpoint. This means that all packets originating from the same logical endpoint that are directed to the same destination node follow the same route across the network, while packets from a different endpoint directed to the same destination node may follow a different path. Figure 6 shows packet routing in an example network. The benefits of this approach is that packet traffic can be distributed across multiple links, while maintaining the order of all packets from the same endpoint. If packets from the same endpoint are allowed to take different paths, it would require a completion buffer which may be expensive in an embedded system. For simplicity, the BlueDBM network does not implement a discovery protocol, and relies on a network configuration file to populate the routing tables.

In order to maintain extremely low network latency, each endpoint is given a choice whether to use end-to-end flow control. If the developer is sure that a particular virtual link will always drain on the receiving end, flow end-to-end flow control can be omitted for that endpoint. However, if the

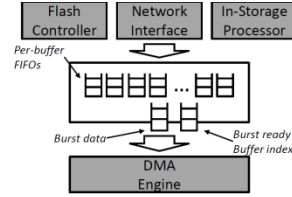


Figure 7: Host-FPGA Interface Over PCIe

receiver fails to drain data for a long time, the link-level back pressure may cause related parts of the network to block. On the other hand, an endpoint can be configured to only send data when there is space on the destination endpoint, which will assure safety but result in higher latency due to flow control packets, and more memory usage for buffers.

3.3. Host Interface

The in-store processing core can be accessed from the host server over either a direct interface that supports RPC and DMA operations, or a file system abstraction built on top of the direct interface. The file system interface is described in detail in Section 4.

In order to parallelize requests and maintain high performance, the host interface provides the software with 128 page buffers, each for reads and writes. When writing a page, the software will request a free write buffer, copy data to the write buffer, and send a write request over RPC with the physical address of the destination flash page. The buffer will be returned to the free queue when the hardware has finished reading the data from the buffer. When reading a page, the software will request a free read buffer, and send a read request over RPC with the physical address of the source flash page. The software will receive an interrupt with the buffer index when the hardware has finished writing to software memory.

Using DMA to write data to the storage device is straightforward to parallelize, but parallelizing reads is a bit more tricky due to the characteristics of flash storage. When writing to storage, the DMA engine on the hardware will read data from each buffer in order in a contiguous stream. So having enough requests in the request queue is enough to make maximum use of the host-side link bandwidth. However, data reads from flash chips on multiple buses in parallel can arrive interleaved at the DMA engine. Because the DMA engine needs to have enough contiguous data for a DMA burst before issuing a DMA burst, some reordering may be required at the DMA engine. This becomes even trickier when the device is using the integrated network to receive data from remote nodes, where they might all be coming from different buses. To fix this issue, we provide dual-ported buffer in hardware which has the semantics of a vector of FIFOs, so that data for each request can be enqueued into its own FIFO until there is enough data for a burst. Figure 7 describes the structure of the host interface for flash reads.

4. Software Interface

In BlueDBM, we aim to provide a set of software interfaces that support the execution of any existing application as well as modified applications that leverage the in-store processors in the system. Furthermore, software layers in BlueDBM must perform flash management functions since we chose to expose a raw flash interface in hardware for higher efficiency (previously discussed in Section 3.1). The software architecture is shown in Figure 8. Three interfaces are supplied to the user application: (i) a file system interface, (ii) a block device driver interface and (iii) an accelerator interface.

We first discuss the file system. Commercial SSDs incorporate a Flash Translation Layer (FTL) inside the flash device controller to manage flash and maintains a block device view to the operating system. However, common file systems manage blocks in a fashion optimized for hard disks. SSDs use the FTL to emulate block device interfaces for compliance with operating systems, performing logical-to-physical mapping and garbage collection, which require large DRAM and incur lots of extra I/Os. Some file systems have tried to remedy this by refactoring the I/O architecture in order to offload most of the FTL functions into a flash-optimized log-structured file system. A prominent example of this is RFS [26]. Unlike conventional FTL designs where the flash characteristics are hidden from the file system, RFS performs some functionality of an FTL, including logical-to-physical address mapping and garbage collection. This achieves better garbage collection efficiency at much lower memory requirement. The file system interface in BlueDBM is built on the same paradigm.

For compatibility with existing software, BlueDBM also offers a full-fledged FTL implemented in the device driver, similar to Fusion IO’s driver. This allows us to use well-known Linux file systems (e.g., ext2/3/4) as well as database systems (directly running on top of a block device) with BlueDBM.

The BlueDBM software allows developers to easily make use of fast in-storage processing without any efforts to write their own custom interfaces manually. Figure 8 shows how user-level applications access hardware accelerators. In the BlueDBM software stack, user-level applications can query the file system for the physical locations of files on the flash (see (1) in Figure 8). This was made possible because the file system maintains the mapping information. Applications can then provide in-storage processors with a stream of physical addresses (see (2)), so that the in-storage processors can directly read data from flash with very low latency (see (3)). The results are sent to software memory and the user application can be notified (see (4)).

It is worth noting that, in BlueDBM, all the user requests, including both user queries and data, are sent to the hardware directly, bypassing almost all of the operating system’s kernel, except for essential driver modules. This helps us to avoid deep OS kernel stacks that often cause long I/O latencies. It is also very common that multiple instances of a user application may

compete for the same hardware acceleration units. For efficient sharing of hardware resources, BlueDBM runs a scheduler that assigns available hardware-acceleration units to competing user-applications. In our implementation, a simple FIFO-based policy is used for request scheduling.

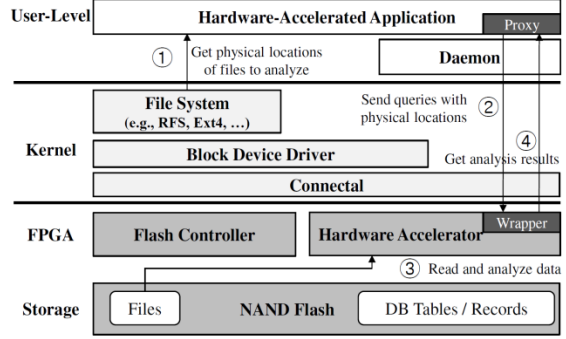


Figure 8: Software Interface

5. Hardware Implementation

We have built a 20-node BlueDBM cluster to explore the capabilities of the architecture. Figure 9 shows the photo of our implementation.

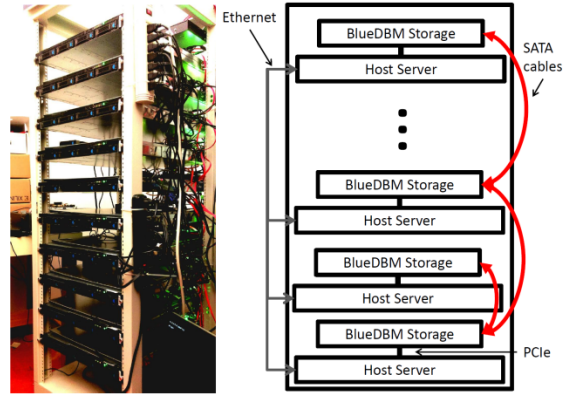


Figure 9: A 20-node BlueDBM Cluster

In our implementation of BlueDBM, we have used a Field Programmable Gate Array (FPGA) to implement the in-store processor and also the flash, host and network controllers. However, the BlueDBM Architecture should not be limited to an FPGA-based implementation. Development of BlueDBM was done in the high-level hardware description language Bluespec. It is possible to develop in-store processors in any hardware description language, as long as they conform to the

interface exposed by the BlueDBM system services. Most of the interfaces are latency-insensitive FIFOs with backpressure. Bluespec provides a lot of support for such interfaces, making in-store accelerator development easier.

The cluster consists of 20 rack-mounted Xeon servers, each with 24 cores and 50GBs of DRAM. Each server also has a Xilinx VC707 FPGA development board connected via a PCIe connection. Each VC707 board hosts two custom-built flash boards with SATA connectors. The VC707 board, coupled with two custom flash boards is mounted on top of each server. The host servers run the Ubuntu distribution of Linux. Figure 10 shows the components of a single node. One of the servers also had a 512GB Samsung M.2 PCIe SSD for performance comparisons.

We used Connectal [24] and its PCIe Gen 1 implementation for the host link. Connectal is a hardware-software codesign framework built by Quanta Research. Connectal reads the interface definition file written by the programmer and generates glue logic between hardware and software. Connectal automatically generates RPC-like interface from developer-provided interface specification, as well as a memory-mapped DMA interface for high bandwidth data transfer. Connectal's PCIe implementation caps our performance at 1.6GB/s reads and 1GB/s writes, which is a reasonable performance for a commodity flash storage device. In the future we will also explore the benefits of a faster host link including later generation PCIe links.

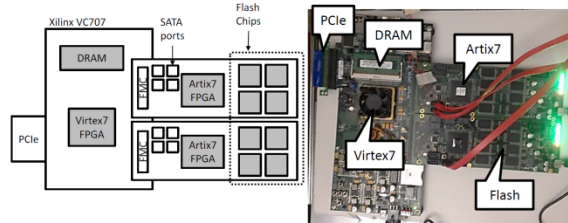


Figure 10: A BlueDBM Storage Node

5.1. Custom Flash Board

We have designed and built a high-capacity custom flash board with high-speed serial connectors, with the help of Quanta Inc., and Xilinx Inc.

Each flash card has 512GBs of NAND flash storage and a Xilinx Artix 7 chip, and plugs into the host FPGA development board via the FPGA Mezzanine Card (FMC) connector. The flash controller and Error Correcting Code (ECC) is implemented on this Artix chip, providing the Virtex 7 FPGA chip on the VC707 a logical error-free access into flash. The communication between the flash board and the Virtex 7 FPGA is done by a 4-lane aurora channel, which is implemented on the GTX/GTP serial transceivers included in each FPGA. This channel can sustain up to 3.3GB/s of bandwidth at $0.5\mu s$

latency. The flash board also hosts 8 SATA connectors, 4 of which pin out the high-speed serial ports on the host Virtex 7 FPGA, and 4 of which pin out the high-speed serial ports on the Artix 7 chip. The serial ports are capable of 10Gbps and 6.6Gbps of bandwidth, respectively.

5.2. Network Infrastructure

In our BlueDBM implementation, the link is implemented over the low-latency serial transceivers. By implementing routing in the hardware and using a very low-latency network fabric, we were able to achieve very high performance, with less than $0.5\mu s$ of latency per network hop, and near 10Gbps of bandwidth per link. Our implementation has a network fan-out of 8 ports per storage node, so the aggregate network bandwidth available to a node reaches up to 8GB/s, including packet overhead.

5.3. Software Interface

Our host interface is implemented using Connectal [24]. Connectal provides a PCIe Gen 1 endpoint and driver pair, and provides up to 1.6GB/s DMA read to host DRAM bandwidth and 1GB/s of DMA write from host DRAM bandwidth. Reading or writing data from the host buffers were done by DMA read/write engines implemented in the Connectal framework. In our BlueDBM implementation, there are four read engines and four write engines each, in order to more easily make maximum use of the PCIe bandwidth.

6. Evaluation

This section evaluates the characteristics of the BlueDBM implementation.

6.1. FPGA Resource Utilization

The FPGA resource usage of each of the two Artix-7 chips are shown in Table 1. 46% of the I/O pins were used either to communicate with the FMC port or to control the flash chips.

Module Name	#	LUTs	Registers	BRAM
Bus Controller	8	7131	4870	21
→ ECC Decoder	2	1790	1233	2
→ Scoreboard	1	1149	780	0
→ PHY	1	1635	607	0
→ ECC Encoder	2	565	222	0
SerDes	1	3061	3463	13
Artix-7 Total		75225 (56%)	62801 (23%)	181 (50%)

Table 1: Flash Controller on Artix 7 Resource Usage

The FPGA resource usage of the Virtex 7 FPGA chip on the VC707 board is shown in Table 2. As it can be seen, there is still enough space for accelerator development on the Virtex FPGA.

Module Name	#	LUTs	Registers	RAMB36	RAMB18
Flash Interface	1	1389	2139	0	0
Network Interface	1	29591	27509	0	0
DRAM Interface	1	11045	7937	0	0
Host Interface	1	88376	46065	169	14
Virtex-7 Total		135271 (45%)	135897 (22%)	224 (22%)	18 (1%)

Table 2: Host Virtex 7 Resource Usage

6.2. Power Consumption

Table 3 shows the overall power consumption of the system, which were estimated using values from the datasheet. Each Xeon server includes 24 cores and 50GBs of DRAM. Thanks to the low power consumption of the FPGA and flash devices, BlueDBM adds less than 20% of power consumption to the system.

Component	Power (Watts)
VC707	30
Flash Board x2	10
Xeon Server	200
Node Total	240

Table 3: BlueDBM Estimated Power Consumption

6.3. Network Performance

We measured the performance of the network by transferring a single stream of 128 bit data packets through multiple nodes across the network in a non-contentious traffic setting. The maximum physical link bandwidth is 10Gbps, and per-hop latency is 0.48 μ s. Figure 11 shows that we are able to sustain 8.2Gbps of bandwidth per stream across multiple network hops. This shows that the protocol overhead is under 18%. The latency is 0.48 μ s per network hop, the end-to-end latency is simply a multiple of network hops to the destination.

Each node in our BlueDBM implementation includes a fan-out of 8 network ports, so each node can have an aggregate full duplex bandwidth of 8.2GB/s. With such a high fan-out, it would be unlikely that a remote node in a rack-class cluster to be over 4 hops, or 2 μ s away. In a naive ring network of 20 nodes with 4 lanes each to next and previous nodes, the average latency to a remote node is 5 hops, or 2.5 μ s. The ring throughput is 32.8 Gbps. *Assuming a flash access latency of 50 μ s, such a network will only add 5% latency in the worst case, giving the illusion of a uniform access storage.*

6.4. Remote Storage Access Latency

We measured the latency of remote storage access by reading an 8K page of data from the following sources using the integrated storage network:

1. ISP-F: From in-store processor to remote flash storage;
2. H-F: From host server to remote flash storage;
3. H-RH-F: From host server to remote flash storage via its host server.

4. H-D: From host server to remote DRAM;

In each case, the request is sent from either the host server or the in-store processor on the local BlueDBM node. In the third and fourth case, the request is processed by the remote server, instead of the remote in-store processor, adding extra latency. However, data is always transferred back via the integrated storage network. We could have also measured the accesses to remote servers via Ethernet, but that latency is at least 100x of the integrated network, and will not be particularly illuminating.

The latency is broken up into four components as shown in Figure 14. First is the local software overhead of accessing the network interface. Second is the storage access latency, or the time it takes for the first data byte to come out of the storage device. Third is the amount of times it takes to transfer the data until the last byte is sent back over the network, and last is the network latency.

Figure 12 shows the exact latency breakdown for each experiment. Notice in all 4 cases, the network latency is insignificant. The data transfer latency is similar except when data is transferred from DRAM (H-D), where it is slightly lower. Notice that except in the case of ISP-F, storage access incurs the additional overhead of PCIe and host software latencies. If we compare ISP-F to H-RH-F, we can see the benefits of an integrated storage network, as the former allows overlapping the latencies of storage and network access.

6.5. Storage Access Bandwidth

We measured the bandwidth of BlueDBM by sending a stream of millions of random read requests for 8KB size pages to local and remote storage nodes, and measuring the elapsed time to process all of the requests. We measured the bandwidth under the following scenarios:

1. Host-Local: Host sends requests to the local flash and all data is streamed returned over PCIe;
2. ISP-Local: Host sends requests to the local flash and all data is consumed at the local in-store processor;
3. ISP-2Nodes: Like ISP-Local except 50% of the requests are sent to a remote flash controller. Only one serial link connects the two nodes;
4. ISP-3Nodes: Like ISP-Local except 33% of the requests are sent to each of the two remote flash controllers. Two serial links connect each remote controller to the local controller.

Figure 13 shows the read bandwidth performance for each of these cases. Our design of the flash card provides 1.2GB/s of bandwidth per card. Therefore in theory, if both cards are kept completely busy 2.4GB/s should be the maximum sustainable bandwidth from the in-store processor, and this is what we observe in the ISP-Local experiment. In our Host-Local experiment, we observed only 1.6GB/s of bandwidth. This is because this is the maximum bandwidth our PCIe implementation can sustain. In ISP-2Nodes, the aggregate bandwidth of two flash devices should add up to 4.8GB/s, but

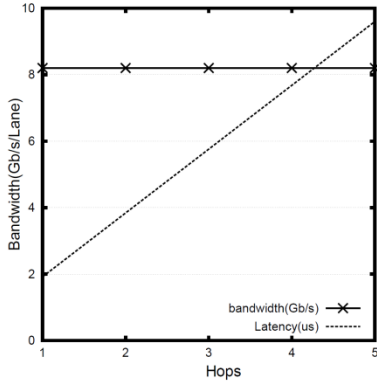


Figure 11: BlueDBM Integrated Network Performance

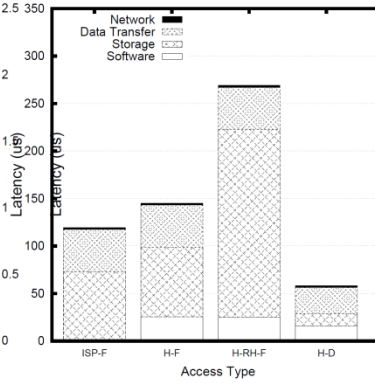


Figure 12: Latency of Remote Data Access in BlueDBM

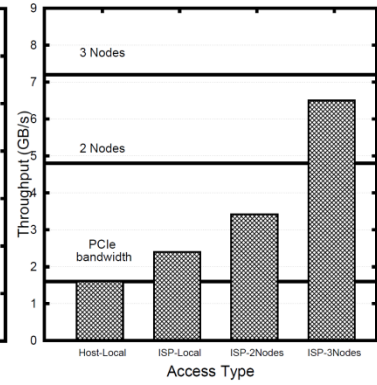


Figure 13: Bandwidth of Data Access in BlueDBM

we only observe about 3.4GB/s, because remote storage access is limited by the single 8Gbps-serial link. In ISP-3Nodes, the aggregate bandwidth of three flash devices should add up to 7.2GB/s, but we only observe about 6.5GB/s because the aggregate bandwidth of the four serial links connecting the remote controllers is limited to 32.8Gbps (≈ 4.1 GB/s).

What these sets of experiments show is that in order to make full use of flash storage, some combination of fast networks, fast host connections and low software overhead is necessary. These requirements can be somewhat mitigated if we make use of in-store computing capabilities, which is what we discuss next.

7. Application Acceleration

In this section, we demonstrate the performance and benefits of the BlueDBM architecture by presenting some accelerator demonstrations.

7.1. Nearest Neighbor Search

Description: Nearest neighbor search is required by many applications, e.g., image querying. One of the modern techniques in this field is Locality Sensitive Hashing [12]. LSH hashes the dataset using multiple hash functions, so that similar data is statistically likely to be hashed to similar buckets. When querying, the query is hashed using the same hash functions, and only the data in the matching buckets are actually

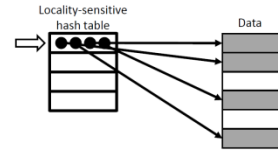


Figure 15: Data Accesses in LSH are Randomly Distributed

compared. The bulk of the work during a query process is traversing hash buckets and reading the corresponding data to perform distance calculation. Because data pointed to by the hash buckets are most likely scattered across the dataset, access patterns are quite random (See Figure 15).

We have built a LSH query accelerator, where all of the data is stored in flash and the distance calculation is done by the in-store processor on the storage device. For simplicity, we assume 8KB data items, and calculate the hamming distance between the query data and each of the items in the hash bucket. The software sends a stream of addresses from a hash bucket along with the query data page, and the system returns the index of the data item most closely matching the query. Since we do not expect any performance difference for queries emanating from two different hash buckets, we simply send out a million nearest-neighbor searches for the same query.

Evaluation: In this study, we were interested in evaluating and comparing the benefits of flash storage (as opposed to DRAM) and in-store processors. We also wanted to compare the BlueDBM design with off-the-shelf SSDs with PCIe interface. The following experiments aim to evaluate the performance of each system during various access patterns, such as random or sequential access, and when accesses are partially serviced by secondary storage.

We have used a commercially available M.2 mPCIe SSD, whose performance, for 8KB accesses, was limited to 600MB/s. Since BlueDBM performance is much higher (2.4GB/s), we also conducted several experiments with BlueDBM throttled to 600MB/s. Since performance should

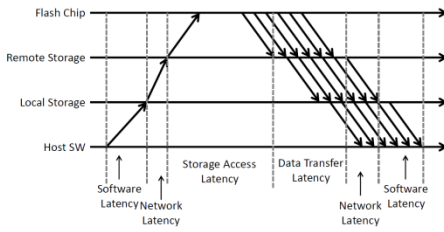


Figure 14: Breakdown of Remote Storage Access Latency

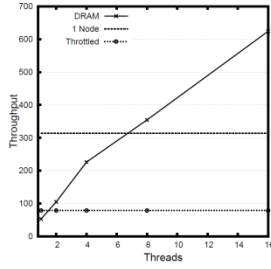


Figure 16: Nearest Neighbor with BlueDBM up to Two Nodes

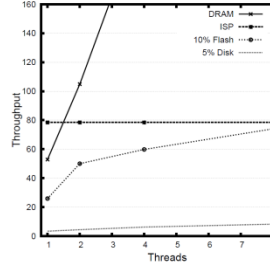


Figure 17: Nearest Neighbor with Mostly DRAM

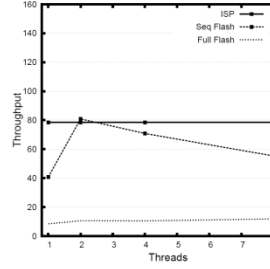


Figure 18: Nearest Neighbor with Off-the-shelf SSD

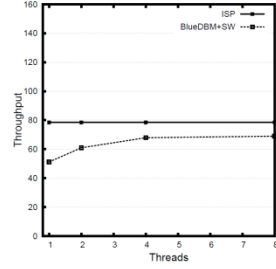


Figure 19: Nearest Neighbor with In-Store Processing

scale linearly with the number of nodes for this application, we concentrated on various configurations in a single node setting:

1. Baseline: BlueDBM with in-store acceleration;
2. Baseline-T: Throttled BlueDBM with in-store acceleration;
3. H-DRAM: Multithread software on multi-core host accessing host DRAM as storage;
4. H-F Throttled: Multithreaded software on multi-core host accessing Throttled BlueDBM as storage;
5. DRAM + 10% Flash: Same as H-DRAM with 10% accesses to SSD;
6. DRAM + 5% Disk: Same as H-DRAM with 5% accesses to HDD;
7. H-RFlash: Multithreaded software on multi-core host accessing Off-the-shelf SSD;
8. H-SFlash: Same as H-RFlash except data accesses are artificially arranged to be sequential.

Figure 16 shows the relative performance of a throttled BlueDBM (Baseline-T) and multithreaded software accessing data on host DRAM (H-DRAM), with Baseline BlueDBM. The baseline performance we observed on BlueDBM was 320K Hamming Comparisons per second. There are two important takeaways from this graph. (1) BlueDBM can keep up with DRAM-resident data for up to 4 threads, because host is getting compute-bound. However, as more threads are added, performance will scale, until DRAM bandwidth becomes the bottleneck. Since DRAM bandwidth as compared to flash bandwidth is very high, DRAM-based processing wins with enough resources. (2) Native flash speed matters i.e., when flash performance is throttled to 1/4th of the maximum, the performance drops accordingly. The relationship between flash performance and application performance will not be so simple if flash was being accessed by software.

To make the comparisons fair, we conducted a set of experiments shown in Figures 17, 18, 19 using throttled BlueDBM as the baseline.

Results of DRAM + 5% Disk and DRAM + 10% Flash experiments shown in Figure 17 show that the performance of ram cloud (H-DRAM) falls off very sharply if even a small fraction of data does not reside in DRAM. Assuming 8 threads,

the performance drops from 350K Hamming Comparisons per second to < 80K and < 10K Hamming Comparisons per second for DRAM + 10% Flash and DRAM + 5% Disk, respectively. At least one commercial vendor has observed similar phenomena and claimed that even when 40% of data fits on DRAM, performance of Hadoop decreases by an order of magnitude [10]. Complex queries on DRAM show high performance only as long as all the data fits in DRAM.

The Off-the-shelf SSD experiment H-RFlash results in Figure 18 showed that its performance is poor as compared to even throttled BlueDBM. However, when we artificially arranged the data accesses to be sequential, the performance improved dramatically, sometimes matching throttled BlueDBM. This suggests that the Off-the-shelf SSD may be optimized for sequential accesses.

Figure 19 comparing Baseline-T and H-F Throttled shows the advantage of accelerators. In this example, the accelerator advantage is at least 20%. Had we not throttled BlueDBM, the advantage would have been 30% or more. This is because while the in-store processor can process data at full flash bandwidth, the software will be bottlenecked by the PCIe bandwidth at 1.6GB/s. We expect this advantage to be larger for applications requiring more complex accelerators. Compared to a fully flash-based execution, BlueDBM performs an order of magnitude faster.

7.2. Graph Traversal

Description: Efficient graph traversal is a very important component of any graph processing system. Fast graph traversal enables solving many problems in graph theory, including maximum flow, shortest path and graph search. It is also a very latency-bound problem because one often cannot predict the next node to visit, until the previous node is visited and processed. We demonstrate the performance benefits of our BlueDBM architecture by implementing distributed graph traversal that takes advantages of the in-store processor and the integrated storage network, which allows extremely low-latency access into both local and remote flash storage.

Evaluation: Graph traversal algorithms often involve dependent lookups. That is, the data from the first request determines

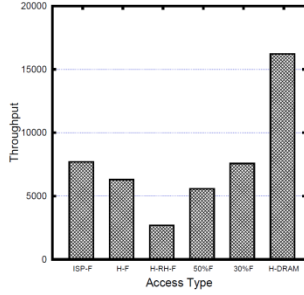


Figure 20: Graph Traversal Performance

the next request, like a linked-list traversal at the page level. Since such traversals are very sensitive to latency, we conducted the experiments with settings that are very similar to the settings in Section 6.4.

1. IPS-F: In-store processor requests data from remote storage over integrated network
2. H-F: Software requests data from remote storage over integrated network
3. H-RH-F: Software requests data from remote software to read from flash
4. DRAM + 50% F: Store requests data from remote software. 50% chance of hitting flash
5. DRAM + 30% F: Store requests data from remote software. 30% chance of hitting flash
6. H-DRAM: Software requests data from remote software. Data read from DRAM

As expected the results in Figure 20 show that the integrated storage network and in-store processor together show almost a factor of 3 performance improvement over generic distributed SSD. This performance difference is large enough that even when 50% of the accesses can be accommodated by DRAM, performance of BlueDBM is still much higher.

The performance difference between *H-F* and *H-RH-F* illustrates the benefits of using the integrated network to reduce a layer of software access. Performance of *ISP-F* compared to *H-F* shows the benefits of further reducing software overhead by having the ISP manage the graph traversal logic.

7.3. String Search

Description: String search is common operation in analytics, often used in database table scans, DNA sequence matching and cheminformatics. It is primarily a sequential read and compare workload. We examine its performance on BlueDBM with assistance from in-store Morris-Pratt (MP) string search engines [31] fully integrated with the file system, flash controller and application software. The software portion of string search initially sets up the accelerator by transferring the target string pattern (needle) and a set of precomputed MP constants over DMA. Then it consults the file system for a list of physical addresses of the files to search (haystack). This list is streamed to the accelerator, which uses these addresses to re-

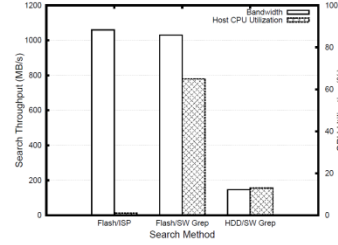


Figure 21: String Search Bandwidth and CPU Utilization

quest for pages from the flash controller. The accelerated MP engines may operate in parallel either by searching multiple files or by dividing up the haystack into equal segments (with some overlaps). This choice depends on the number of files and size of each file. Since 4 read commands can saturate a single flash bus, we use 4 engines per bus to maximize the flash bandwidth. Only search results are returned to the server.

Evaluation: We compared our implementation of hardware-accelerated string search running on BlueDBM to the Linux Grep utility querying for exact string matches running on both SSD and hard disk. Processing bandwidth and server CPU utilizations are shown in Figure 21. We observe that the parallel MP engines in BlueDBM are able to process a search at 1.1GB/s, which is 92% of the maximum sequential bandwidth a single flash board. Using BlueDBM, the query consumes almost no CPU cycles on the host server since the query is entirely offloaded and only the location of matched strings are returned, which we assume is a tiny fraction of the file (0.01% is used in our experiments). This is 7.5x faster than software string search (Grep) on hard disks, which is I/O bound by disk bandwidth and consumes 13% CPU. On SSD, software string search remains I/O bound by the storage device, but CPU utilization increases significantly to 65% even for this type of simple streaming compare operation. This high utilization is problematic because string search is often only a small portion of more complex analytic queries that can quickly become compute bound. As we have shown in the results, BlueDBM can effectively alleviate this by offloading search to the in-store processor thereby freeing up the server CPU for other tasks.

8. Conclusion and Future Work

We have presented BlueDBM, an appliance for Big Data analytics that uses flash storage, in-store processing and integrated networks for cost-effective analytics of large datasets. A rack-size BlueDBM system is likely to be an order of magnitude cheaper and less power hungry than a cloud based system with enough DRAM to accommodate 10TB to 20TB of data. We have demonstrated the performance benefits of BlueDBM using simple examples on large amounts of data in comparison to a generic flash-based system without such architectural improvements. We have also shown that the performance of

a system which relies on data being resident in DRAM, falls rapidly if even a small fraction of data has to reside in secondary storage. BlueDBM like architecture does not suffer from this problem because flash based systems with 10TB to 20TB of storage are very affordable.

Our current implementation uses an FPGA to implement most of the new architectural features, that is, in-store processors, integrated network routers, flash controllers. It is straightforward to implement most of these features using ASICs and provide some in-store computing capability via general-purpose processors. This will simultaneously improve the performance and lower the power consumption even further. Notwithstanding such developments we are developing tools to make it easy to develop in-store processors for the reconfigurable logic inside BlueDBM.

We are currently developing or planning to develop several new applications including: *SQL Database Acceleration* by offloading query processing and filtering to in-store processors, *Sparse-Matrix Based Linear Algebra Acceleration* and *BlueDBM-Optimized MapReduce*, which attempts to optimize data flow of MapReduce to best fit an SSD-based cluster with in-store processors. We plan to collaborate with other research groups to explore more applications.

9. Acknowledgements

This work was partially funded by Quanta (Agmt. Dtd. 04/01/05), Samsung (Res. Agmt. Eff. 01/01/12), Lincoln Laboratory (PO7000261350), and Intel (Agmt. Eff. 07/23/12). We also thank Xilinx for their generous donation of VC707 FPGA boards and FPGA design expertise.

References

- [1] N. Agrawal, V. Prabhakaran, T. Wobber, J. D. Davis, M. Manasse, and R. Panigrahy, "Design tradeoffs for ssd performance," in *USENIX 2008 Annual Technical Conference on Annual Technical Conference*, ser. ATC'08. Berkeley, CA, USA: USENIX Association, 2008, pp. 57–70. Available: <http://dl.acm.org/citation.cfm?id=1404014.1404019>
- [2] I. T. Association, *Infiniband*, 2014 (Accessed November 18, 2014). Available: <http://www.infinibanda.org>
- [3] J. Banerjee, D. Hsiao, and K. Kannan, "Dbc: A database computer for very large databases," *Computers, IEEE Transactions on*, vol. C-28, no. 6, pp. 414–429, June 1979.
- [4] A. M. Caulfield, A. De, J. Coburn, T. I. Mollow, R. K. Gupta, and S. Swanson, "Moneta: A high-performance storage array architecture for next-generation, non-volatile memories," in *Proceedings of the 2010 43rd Annual IEEE/ACM International Symposium on Microarchitecture*, ser. MICRO '43. Washington, DC, USA: IEEE Computer Society, 2010, pp. 385–395. Available: <http://dx.doi.org/10.1109/MICRO.2010.33>
- [5] A. M. Caulfield and S. Swanson, "Quicksan: A storage area network for fast, distributed, solid state disks," *SIGARCH Comput. Archit. News*, vol. 41, no. 3, pp. 464–474, Jun. 2013. Available: <http://doi.acm.org/10.1145/2508148.2485962>
- [6] B. Y. Cho, W. S. Jeong, D. Oh, and W. W. Ro, "Xsd: Accelerating mapreduce by harnessing the gpu inside an ssd," 2013.
- [7] E. S. Chung, J. D. Davis, and J. Lee, "Linqits: Big data on little clients," in *Proceedings of the 40th Annual International Symposium on Computer Architecture*, ser. ISCA '13. New York, NY, USA: ACM, 2013, pp. 261–272. Available: <http://doi.acm.org/10.1145/2485922.2485945>
- [8] J. Dai, "Toward efficient provisioning and performance tuning for hadoop," *Proceedings of the Apache Asia Roadshow*, vol. 2010, pp. 14–15, 2010.
- [9] J. Do, Y.-S. Kee, J. M. Patel, C. Park, K. Park, and D. J. DeWitt, "Query processing on smart ssds: Opportunities and challenges," in *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD '13. New York, NY, USA: ACM, 2013, pp. 1221–1230. Available: <http://doi.acm.org/10.1145/2463676.2465295>
- [10] FusionIO, *Using HBase with ioMemory*, 2012 (Accessed November 22, 2014). Available: <http://www.fusionio.com/white-papers/using-hbase-with-iomemory>
- [11] FusionIO, *FusionIO*, 2014 (Accessed November 18, 2014). Available: <http://www.fusionio.com>
- [12] A. Gionis, P. Indyk, R. Motwani *et al.*, "Similarity search in high dimensions via hashing," in *Vldb*, vol. 99, 1999, pp. 518–529.
- [13] Google, *Google Flu Trends*, 2011 (Accessed November 18, 2014). Available: <http://www.google.org/flu-trends>
- [14] S. Hardock, I. Petrov, R. Gottstein, and A. Buchmann, "Noftl: Database systems on ftl-less flash storage," *Proc. VLDB Endow.*, vol. 6, no. 12, pp. 1278–1281, Aug. 2013. Available: <http://dx.doi.org/10.14778/2536274.2536295>
- [15] M. Honda, F. Huici, C. Raiciu, J. Araujo, and L. Rizzo, "Rekindling network protocol innovation with user-level stacks," *SIGCOMM Comput. Commun. Rev.*, vol. 44, no. 2, pp. 52–58, Apr. 2014. Available: <http://doi.acm.org/10.1145/2602204.2602212>
- [16] Intel, *Intel Solid-State Drive Data Center Family for PCIe*, 2014 (Accessed November 18, 2014). Available: <http://www.intel.com/content/www/us/en/solid-state-drives/intel-ssd-dc-family-for-pcie.html>
- [17] N. S. Islam, M. W. Rahman, J. Jose, R. Rajachandrasekar, H. Wang, H. Subramoni, C. Murthy, and D. K. Panda, "High performance rdma-based design of hdfs over infiniband," in *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, ser. SC '12. Los Alamitos, CA, USA: IEEE Computer Society Press, 2012, pp. 35:1–35:35. Available: <http://dl.acm.org/citation.cfm?id=2388996.2389044>
- [18] Z. István, L. Woods, and G. Alonso, "Histograms as a side effect of data movement for big data," in *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*. ACM, 2014, pp. 1567–1578.
- [19] E. Y. Jeong, S. Woo, M. Jamshed, H. Jeong, S. Ihm, D. Han, and K. Park, "mtep: A highly scalable user-level tcp stack for multicore systems," in *Proceedings of the 11th USENIX Conference on Networked Systems Design and Implementation*, ser. NSDI'14. Berkeley, CA, USA: USENIX Association, 2014, pp. 489–502. Available: <http://dl.acm.org/citation.cfm?id=2616448.2616493>

- [20] S.-W. Jun, M. Liu, K. E. Fleming, and Arvind, "Scalable multi-access flash store for big data analytics," in *Proceedings of the 2014 ACM/SIGDA International Symposium on Field-programmable Gate Arrays*, ser. FPGA '14. New York, NY, USA: ACM, 2014, pp. 55–64. Available: <http://doi.acm.org/10.1145/2554688.2554789>
- [21] S.-H. Kang, D.-H. Koo, W.-H. Kang, and S.-W. Lee, "A case for flash memory ssd in hadoop applications," *International Journal of Control and Automation*, vol. 6, no. 1, 2013.
- [22] Y. Kang, Y.-s. Kee, E. L. Miller, and C. Park, "Enabling cost-effective data processing with smart ssd," in *Mass Storage Systems and Technologies (MSST), 2013 IEEE 29th Symposium on*. IEEE, 2013, pp. 1–12.
- [23] K. Keeton, D. A. Patterson, and J. M. Hellerstein, "A case for intelligent disks (idisks)," *SIGMOD Rec.*, vol. 27, no. 3, pp. 42–52, Sep. 1998. Available: <http://doi.acm.org/10.1145/290593.290602>
- [24] M. King, J. Hicks, and J. Ankcorn, "Software-driven hardware development," in *Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, ser. FPGA '15. New York, NY, USA: ACM, 2015, pp. 13–22. Available: <http://doi.acm.org/10.1145/2684746.2689064>
- [25] O. Kocherber, B. Grot, J. Picorel, B. Falsafi, K. Lim, and P. Ranganathan, "Meet the walkers: Accelerating index traversals for in-memory databases," in *Proceedings of the 46th Annual IEEE/ACM International Symposium on Microarchitecture*, ser. MICRO-46. New York, NY, USA: ACM, 2013, pp. 468–479. Available: <http://doi.acm.org/10.1145/2540708.2540748>
- [26] S. Lee, J. Kim, and Arvind, "Refactored design of i/o architecture for flash storage," *Computer Architecture Letters*, vol. PP, no. 99, pp. 1–1, 2014.
- [27] S.-W. Lee, B. Moon, C. Park, J.-M. Kim, and S.-W. Kim, "A case for flash memory ssd in enterprise database applications," in *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*. ACM, 2008, pp. 1075–1086.
- [28] H. Leilich, G. Stiege, and H. C. Zeidler, "A search processor for data base management systems," in *Fourth International Conference on Very Large Data Bases, September 13-15, 1978, West Berlin, Germany*, 1978, pp. 280–287.
- [29] J. Liu, J. Wu, S. P. Kini, P. Wyckoff, and D. K. Panda, "High performance rdma-based mpi implementation over infiniband," in *Proceedings of the 17th Annual International Conference on Supercomputing*, ser. ICS '03. New York, NY, USA: ACM, 2003, pp. 295–304. Available: <http://doi.acm.org/10.1145/782814.782855>
- [30] V. Memory, *Violin Memory*, 2014 (Accessed November 18, 2014). Available: <http://www.violin-memory.com>
- [31] J. Morris Jr and V. Pratt, *A linear pattern-matching algorithm*, 1970.
- [32] R. Mueller, J. Teubner, and G. Alonso, "Streams on wires: A query compiler for fpgas," *Proc. VLDB Endow.*, vol. 2, no. 1, pp. 229–240, Aug. 2009. Available: <http://dx.doi.org/10.14778/1687627.1687654>
- [33] Oracle, *Exadata Database Machine*, 2014 (Accessed November 18, 2014). Available: <https://www.oracle.com/engineered-systems/exadata/index.html>
- [34] J. Ousterhout, P. Agrawal, D. Erickson, C. Kozyrakis, J. Leverich, D. Mazières, S. Mitra, A. Narayanan, G. Parulkar, M. Rosenblum, S. M. Rumble, E. Stratmann, and R. Stutsman, "The case for ramclouds: Scalable high-performance storage entirely in dram," *SIGOPS Oper. Syst. Rev.*, vol. 43, no. 4, pp. 92–105, Jan. 2010. Available: <http://doi.acm.org/10.1145/1713254.1713276>
- [35] E. A. Ozkaran, S. A. Schuster, and K. C. Smith, "RAP - an associative processor for database management," in *American Federation of Information Processing Societies: 1975 National Computer Conference, 19-22 May 1975, Anaheim, CA, USA*, 1975, pp. 379–387. Available: <http://doi.acm.org/10.1145/1499949.1500024>
- [36] A. Putnam, A. M. Caulfield, E. S. Chung, D. Chiou, K. Constantinides, J. Demme, H. Esmaeilzadeh, J. Fowers, G. Prashanth, G. Jan, G. Michael, H. S. Hauck, S. Heil, A. Hormati, J.-Y. Kim, S. Lanka, J. Larus, E. Peterson, S. Pope, A. Smith, J. Thong, P. Yi, and X. D. Burger, "A reconfigurable fabric for accelerating large-scale datacenter services," *SIGARCH Comput. Archit. News*, vol. 42, no. 3, pp. 13–24, Jun. 2014. Available: <http://doi.acm.org/10.1145/2678373.2665678>
- [37] M. W.-u. Rahman, N. Islam, X. Lu, J. Jose, H. Subramoni, H. Wang, and D. Panda, "High-performance rdma-based design of hadoop mapreduce over infiniband," in *International Workshop on High Performance Data Intensive Computing (HPDIC), in conjunction with IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2013.
- [38] M. W.-u. Rahman, X. Lu, N. S. Islam, and D. K. D. Panda, "Homr: A hybrid approach to exploit maximum overlapping in mapreduce over high performance interconnects," in *Proceedings of the 28th ACM International Conference on Supercomputing*, ser. ICS '14. New York, NY, USA: ACM, 2014, pp. 33–42. Available: <http://doi.acm.org/10.1145/2597652.2597684>
- [39] S. M. Rumble, A. Kejriwal, and J. Ousterhout, "Log-structured memory for dram-based storage," in *Proceedings of the 12th USENIX Conference on File and Storage Technologies*, ser. FAST '14. Berkeley, CA, USA: USENIX Association, 2014, pp. 1–16. Available: <http://dl.acm.org/citation.cfm?id=2591305.2591307>
- [40] SanDisk, *Sandisk ZetaScale Software*, 2014 (Accessed November 22, 2014). Available: <http://www.sandisk.com/enterprise/zetascale/>
- [41] S. Seshadri, M. Gahagan, S. Bhaskaran, T. Bunker, A. De, Y. Jin, Y. Liu, and S. Swanson, "Willow: A user-programmable ssd," in *Proceedings of the 11th USENIX Conference on Operating Systems Design and Implementation*, ser. OSDI '14. Berkeley, CA, USA: USENIX Association, 2014, pp. 67–80. Available: <http://dl.acm.org/citation.cfm?id=2685048.2685055>
- [42] M. Singh and B. Leonhardi, "Introduction to the ibm netezza warehouse appliance," in *Proceedings of the 2011 Conference of the Center for Advanced Studies on Collaborative Research*, ser. CASCON '11. Riverton, NJ, USA: IBM Corp., 2011, pp. 385–386. Available: <http://dl.acm.org/citation.cfm?id=2093889.2093965>
- [43] J. Spiegel, M. McKenna, G. Lakshman, and P. Nordstrom, "Method and system for anticipatory package shipping," Dec. 27 2011, uS Patent 8,086,546. Available: <http://www.google.com/patents/US8086546>
- [44] B. Sukhwani, H. Min, M. Thoennes, P. Dube, B. Iyer, B. Brezzo, D. Dillenberger, and S. Asaad, "Database analytics acceleration using fpgas," in *Proceedings of the 21st International Conference on Parallel Architectures and Compilation Techniques*, ser. PACT '12. New York, NY, USA: ACM, 2012, pp. 411–420. Available: <http://doi.acm.org/10.1145/2370816.2370874>
- [45] D. Technologies, *Diablo Technologies*, 2014 (Accessed November 18, 2014). Available: <http://www.diablo-technologies.com/>
- [46] T. S. Woodall, G. M. Shipman, G. Bosilca, R. L. Graham, and A. B. Maccabe, "High performance rdma protocols in hpc," in *Proceedings of the 13th European PVM/MPI User's Group Conference on Recent Advances in Parallel Virtual Machine and Message Passing Interface*, ser. EuroPVM/MPI '06. Berlin, Heidelberg: Springer-Verlag, 2006, pp. 76–85. Available: http://dx.doi.org/10.1007/11846802_18
- [47] L. Woods, Z. Istvan, and G. Alonso, "Hybrid fpga-accelerated sql query processing," in *Field Programmable Logic and Applications (FPL), 2013 23rd International Conference on*, Sept 2013, pp. 1–1.
- [48] L. Woods, Z. Istvan, and G. Alonso, "Ibex - an intelligent storage engine with support for advanced sql off-loading," in *Proceedings of the 40th International Conference on Very Large Data Bases*, ser. VLDB '14, 2014, pp. 963–974.
- [49] L. Wu, A. Lottarini, T. K. Paine, M. A. Kim, and K. A. Ross, "Q100: The architecture and design of a database processing unit," in *Proceedings of the 19th International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS '14. New York, NY, USA: ACM, 2014, pp. 255–268. Available: <http://doi.acm.org/10.1145/2541940.2541961>

This page intentionally left blank.

APPENDIX B

A Transport-Layer Network for Distributed FPGA Platforms

Sang-Woo Jun, Ming Liu, Shuotao Xu, Arvind
Department of Electrical Engineering and Computer Science
Massachusetts Institute of Technology
{wjun, ml, shuotao, arvind}@csail.mit.edu

Abstract—We present a transport-layer network that aids developers in building safe, high-performance distributed FPGA applications. Two essential features of such a network are virtual channels and end-to-end flow control. Our network implements these features, taking advantage of the low error characteristic of a rack level FPGA network to implement a low overhead credit based end-to-end flow control. Our design has many parameters in the source code which can be set at the time of FPGA synthesis, to provide flexibility in setting buffer size and flow control credits to make best use of scarce on-chip memory resources and match the traffic pattern of a virtual channel. Our prototype cluster, which is composed of 20 Xilinx VC707 boards, each with 4 20Gb/s serial links, achieves effective bandwidth of 85% of the maximum physical bandwidth, and a latency of 0.5us per hop. User feedback suggest that these features make distributed application development significantly easier.

I. INTRODUCTION

In order to tackle large data intensive applications, many modern FPGA-based deployments are exploring the use of FPGA clusters, where a network of FPGAs are deployed and a large body of work is distributed across the FPGAs. A network protocol for an FPGA cluster largely has three important criteria: (1) it must be easily usable by an application developer, (2) It must have high performance with low latency, and (3) it must consume only a small amount of scarce on-chip FPGA memory.

Two essential features for a usable network implementation are virtual channels and end-to-end flow control, corresponding to the transport layer of the OSI network model. Without these features, the developer would have to manually manage channel multiplexing and deadlock management, making the development of high performance distributed FPGA applications difficult.

Due to the high engineering and performance overhead of existing network solutions, many inter-FPGA networks on a distributed FPGA deployment are implemented using low-overhead link-layer protocols such as Aurora using multi-gigabit serial transceivers included in the FPGA. Many existing network implementations using this network fabric often provide link and network level interfaces, but they rarely provide higher-level functionality such as end-to-end flow control.

For deadlock-free operations, all virtual channels need separate packet buffers which are large enough to mask network latency as well as bursts from multiple sources. Scarcity of on-chip memory resources prevent safe over-allocation of packet buffers, and using off-chip DRAM also consumes a lot of precious DRAM bandwidth. A solution may be clever allocation

of buffer space by allowing different amount of buffers for different channels. Application developers can adjust the buffer space per channel to meet the performance criteria without increasing the total buffer requirement.

This paper presents the design and implementation results of a transport level network for a cluster of FPGAs. Our transport layer is parameterized such that flow control features for each virtual channel can be configured at FPGA synthesis time. Parameters include buffer size and flow control credits. We demonstrate that a parameterized transport-layer implementation can achieve high performance in a distributed FPGA environment while maintaining a small BRAM footprint, by adjusting a few parameters to best fit the usage characteristics of a virtual channel. In our router, we make use of the high reliability of the serial link and deterministic routing to ensure lossless in-order arriving of packets, greatly simplifying the transport layer protocol.

We have implemented a prototype of our network on a cluster of 20 Xilinx VC707 FPGA development boards, with 4 20Gb/s serial links each. Our prototype achieves an effective bandwidth of 17Gb/s per link, which is 85% of maximum physical link bandwidth, at a latency of 0.5us.

The rest of the paper is organized as follows: Section II covers the previous and related work. Section III describes our implementation of the network and transport layer, and Section IV describes the details of a prototype implementation of the network. Section V presents the performance evaluation of our implementation, and conclude in Section VI.

II. RELATED WORK

FPGAs offer very desirable performance and power characteristics, but modern data-intensive applications often require more resources that are available on a single FPGA chip. As a result, exploration of distributed FPGA computing systems is gaining popularity. The scale of distributed FPGA deployments range from a cluster-in-a-box systems such as BlueHive [1], to rack-level deployments such as Maxwell [2], to datacenter scale deployments such as Catapult [3]. Some have also attempted heterogeneous deployments including GPUs and FPGAs [4], or to insert FPGA accelerators into the storage datapath [5], [6]. Such systems offer a much better power performance characteristics over their off-the-shelf server counterparts.

The TCP/IP network protocol stack is by far the most popular protocol for internetworking computer systems, but it may not be a good fit for inter-FPGA communication as

it is a complex and resource-heavy protocol designed for an unpredictable network such as the internet. Some FPGA cluster projects have used Ethernet's physical and data link layers for its network, but full implementation of the TCP/IP stack is rare unless it has to interact with a legacy interface [7]. Datacenter scale protocols such as Infiniband [8] implement a more efficient transport layer protocol on top of a more reliable network layer implementation. It also offloads major parts of the protocol to the NIC to achieve higher performance. Some have modified the TCP protocol [9] to significantly reduce the packet buffer size using intelligent congestion control.

BlueLink [10] demonstrated that a new protocol using high-speed serial links has a better area-performance characteristics than implementing existing network protocols. Many distributed FPGA computing systems have demonstrated high performance with FPGA nodes networked over such high-speed serial links [11], [2]. Some have developed meta language compilers that generate application-specific network logic with features such as flow control from separate network specifications [12].

III. NETWORK ARCHITECTURE

The overall architecture of the network components can be seen in Figure 1. The network architecture can be divided largely into two parts, the network layer and the transport layer. The network layer is implemented in the form of the router, and the transport layer is implemented in the endpoints that are chained to the router interface. Flow control is implemented in both layers.

The distributed application components communicate with remote nodes using the network *endpoints*. Endpoints expose send and receive interfaces, and behaves like a FIFO, in that it blocks when it cannot safely send any more packets. Many endpoints can be instantiated, resources permitting, and each endpoint can have a different *type*, meaning it can expose send and receive interfaces of different bit widths.

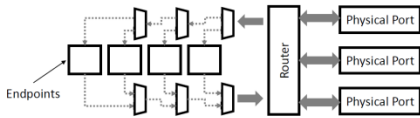


Fig. 1: Network Architecture

A. Network Layer

The network layer implements lossless, in-order packet routing, which assures that packets always arrive in the order they were sent. This removes the need for try-resend or reordering functionalities at the transport layer, allowing a much simpler design.

The router implements in-order routing by being deterministic, in that a packet from a certain endpoint of a certain source node being delivered to a certain destination node will always travel through the same path. Parallelism is achieved by distributing packets from different endpoints to different paths. The router is oblivious to the existence of multiple network endpoints or virtual channels they represent. The

router only deals with routing individual packets, and higher level functions such as virtual channels and end-to-end flow control is implemented by the endpoints, which are organized into a chain to reduce fan-in on the FPGA.

A packet consists of four fields: source node ID, destination node ID, endpoint ID and payload data. Destination node ID and payload data is supplied by the user. Source node ID is supplied by the router, and the endpoint ID is determined by the position of an endpoint in the endpoint chain. Endpoint ID is filled out by the endpoint chain when a packet is injected into it, and it is used to direct a packet to the correct endpoint at the receiving side.

B. Transport Layer

Virtual channels multiplex a single physical network link to provide the logical interface of multiple links. Figure 2 describes the flow of packets in such an environment. Our network implements a per-channel end-to-end flow control, so that a sender can only send data onto the network when it is guaranteed that the receiving endpoint has enough buffer space to accommodate it. The transport layer is implemented in individual endpoints, and its design aims to provide a very low latency and efficient memory space usage. Each design can have multiple instantiations of endpoints, parameterized differently.

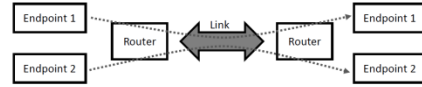


Fig. 2: Packet Flow in Virtual Channels

The structure of an endpoint is described in Figure 3. Whenever a packet is received by an endpoint, it checks a table of packets received per source node to determine if it is time to send a flow control credit to the source node. If the send budget of the source node is predicted to have become small enough and there is enough space on the local receive buffer, it enters a packet into the ack queue and marks the amount of space as allocated.

In order to maintain maximum bandwidth, flow control packets must be received before the send budget of the source node runs out. However, it is often not possible to provide a large enough buffer to conservatively accommodate the flow control packet's round trip latency for all nodes in the system.

Under such constraints, each endpoint can have a different flow control configuration that attempts to best suit its usage. For example, for some endpoints the expected traffic pattern may be that most of the data transfer happens between a pair of two nodes. In such a case, it might be effective to have a very low granularity flow control, so that a large buffer is allocated on request, but the total receive buffer may be small. On the other hand, if many nodes are expected to send data to one node, a fine granularity flow control and a large buffer may be required for performance. If the endpoint is used for low-bandwidth traffic such as commands, the buffer size and granularity can be set to a small value. To enable such control, endpoints are initialized with the parameters described in Table I.

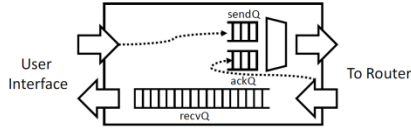


Fig. 3: Endpoint Architecture

Initially, all nodes start with a small send budget ($initBudget$) to all remote nodes, and therefore the actual size of the receive packet buffer is $initBudget \times nodeCount$ slots larger than the $BufferSize$ parameter. When the first packet arrives, space in the receive buffer is allocated to the source node and a flow control packet is sent. The endpoint can choose to periodically allocate only $initBudget$ size buffers, instead of $FlowCredit$ in an attempt to more fairly allocate buffer space across source nodes. Yielding buffers like this achieves better buffer usage when many nodes are going to send data to one node.

The network infrastructure also provides an unmanaged endpoint which does not implement a transport layer protocol. The unmanaged endpoint can sustain the highest bandwidth and lowest latency as it does not check if the receiving endpoint has available buffers before sending packets. It should be used very carefully, since if the arriving data is not always immediately consumed and dequeued from the receiving buffer, it may cause the entire network to block.

IV. IMPLEMENTATION DETAILS

We have implemented a prototype of the network described using a cluster of machines. Each node in the cluster consists of one Intel Xeon-based server, Xilinx VC707 FPGA development board, and a network expansion card which pinned out eight GTX multi-gigabit serial transceivers. Two lanes were grouped together to form a channel, resulting in a fan-out of up to 4 channels per node. A node can implement any direct network with a fan-out per node of 4 or less, such as a 2D mesh or torus.

The link latency of an aurora link based on the GTX multi-gigabit transceiver was measured to be around 0.48us, which translates to about 75 cycles on the 6.4ns user clock.

A. FPGA Resource Utilization

We have measured the FPGA resource utilization of our network using a simple setup with two endpoints: one high speed endpoint with larger flow control credit and buffer size (Credit size of 200 and buffer size of 1024 packets), and one small endpoint with smaller buffers. The endpoint row in the table below described the larger endpoint. The router component includes the chaining logic used to link the endpoints to it. The user logic was clocked at 125MHz.

Parameter	Description
BufferSize	Size of the total allocated buffer space
FlowOffset	Offset of flow control packet transmission
FlowCredit	Number of packets each flow control credit represents

TABLE I: Endpoint Parameters

Component	LUTs	RAMB36
Aurora Link	4843	36
Router	3743	0
Endpoint ($\times 2$)	753	3
Total	10092	42
Virtex 7 Percentage	(3%)	(4%)

TABLE II: FPGA Resource Utilization

V. PERFORMANCE EVALUATION

We demonstrate that the performance of the network does not suffer from the addition of transport-layer network functions. We measured the bandwidth and latency of the network under various configurations, and show that our network can usually achieve a bandwidth of 17Gb/s, which is 85% of the maximum physical link bandwidth. This performance is reasonable considering the packet header and flow control overhead.

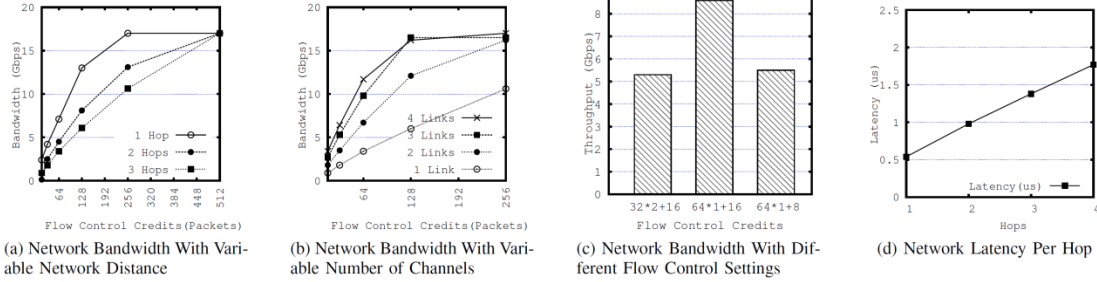
Single Endpoint Over Multiple Hops : We measured the bandwidth of the network implementation by measuring the time it takes for a single endpoint to send a large amount of data to remote nodes variable hops away, under various flow control settings. Larger credit settings mean a larger buffer size is required. Flow control offset was set to be half of the credit size. The results can be seen in Figure 4a.

When the flow control credit was small, performance of the network was lower when going over a longer network distance. This is because the round trip latency over multiple hops is longer than the time it takes to deplete the send buffer, resulting in idle cycles when no data can safely be sent over the network. With the low network latency of the serial links, maximum bandwidth over 3 network hops could be achieved using a single endpoint when the flow control credit is over 512 packets large.

Multiple Endpoints Over Multiple Hops : Since most interesting distributed FPGA applications will have more than one network endpoint, maximum network performance can be achieved even when a single endpoint's flow control credit setting is large enough. We measured the aggregate network bandwidth of a varying number of endpoints sending data to a node three network hops away. We also measured the performance with varying flow control credit sizes. Flow control offset was set to be half of the credit size. The results can be seen in Figure 4b. It shows that a collection of smaller sized endpoints can saturate the network by filling in each others' idle cycles.

Buffer Size and Flow Control Offset : Endpoints can be characterized not only by its flow control credit size, but also by the flow control offset and buffer size parameters. The same amount of buffer space can also be allocated to a different number of nodes under different flow control credit settings. Setting a smaller offset has the risk of incurring idle time by delivering a flow control packet too late, but a large offset requires a larger buffer to accommodate earlier buffer allocation.

We measured the effect of such parameters by having three nodes send a stream of packets to the same remote node. We tested three scenarios, described in Table III. Two had the same



total buffer size organized into different organizations, and one had a smaller buffer. In the first scenario, the three source nodes will be contending to be scheduled into the two possible slots, where in the latter two scenarios they will be contending for one 64 packet slot.

Setting	Description
32*2+16	Buffer has space for two 32 packet blocks, with offset of 16
64*1+16	Buffer has space for one 64 packet block, with offset of 16
64*1+8	Buffer has space for one 64 packet block, with offset of 8

TABLE III: Flow Control Parameters

The results can be seen in Figure 4c. It shows that even with the same buffer size, having a larger credit setting is beneficial to a small buffer configuration. The difference is pronounced enough that even reducing buffer usage further by making the offset smaller results in a better performance compared to the configuration with smaller credit sizes.

Multi-Hop Latency : Network latency was measured by measuring the round-trip latency by sending a packet to nodes of varying distances, where the user logic immediately sends the packet back to the original sender. The results can be seen in Figure 4d. We show a consistent latency of less than 0.5us per hop.

VI. CONCLUSION

In this paper, we have presented our design of a parameterized, low overhead transport-layer network that provides useful features such as virtual channels and end-to-end flow control. Our network takes advantage of the high reliability of the high-speed serial links, which are integrated in the FPGA fabric, to implement a lossless in-order network layer, which allowed us to simplify the transport layer and use less FPGA resources. The design of the transport layer is parameterized, so that the developer can choose to use less resources while meeting the performance requirements of the individual endpoint. Our prototype implementation demonstrated a high performance in an FPGA cluster setting. We predict that our network will accelerate future research of distributed FPGA applications.

ACKNOWLEDGEMENT

This work was partially funded by Quanta (Agmt. Dtd. 04/01/05) and Lincoln Laboratory(PO7000261350). We also thank Xilinx for their generous donation of VC707 FPGA boards and FPGA design expertise.

REFERENCES

- [1] S. Moore, P. Fox, S. Marsh, A. Markettos, and A. Mujumdar, "Bluehive - a field-programable custom computing machine for extreme-scale real-time neural network simulation," in *Field-Programmable Custom Computing Machines (FCCM)*, 2012 IEEE 20th Annual International Symposium on, April 2012, pp. 133–140.
- [2] R. Baxter, S. Booth, M. Bull, G. Cawood, J. Perry, M. Parsons, A. Simpson, A. Trew, A. McCormick, G. Smart, R. Smart, A. Cantle, R. Chamberlain, and G. Genest, "Maxwell - a 64 fpga supercomputer," in *Adaptive Hardware and Systems, 2007. AHS 2007. Second NASA/ESA Conference on*, Aug 2007, pp. 287–294.
- [3] A. Putnam, A. M. Caulfield, E. S. Chung, D. Chiou, K. Constantinides, J. Demme, H. Esmailzadeh, J. Fowers, G. Prashanth, G. Jan, G. Michael, H. S. Hauck, S. Heil, A. Hormati, J.-Y. Kim, S. Lanka, J. Larus, E. Peterson, S. Pope, A. Smith, J. Thong, P. Yi, and X. D. Burger, "A reconfigurable fabric for accelerating large-scale datacenter services," *SIGARCH Comput. Archit. News*, vol. 42, no. 3, pp. 13–24, Jun. 2014.
- [4] K. H. Tsoi and W. Luk, "Axel: A heterogeneous cluster with fpgas and gpus," in *Proceedings of the 18th Annual ACM/SIGDA International Symposium on Field Programmable Gate Arrays*, ser. FPGA '10. New York, NY, USA: ACM, 2010, pp. 115–124.
- [5] S.-W. Jun, M. Liu, K. E. Fleming, and Arvind, "Scalable multi-access flash store for big data analytics," in *Proceedings of the 2014 ACM/SIGDA International Symposium on Field-programmable Gate Arrays*, ser. FPGA '14. New York, NY, USA: ACM, 2014, pp. 55–64.
- [6] S.-W. Jun, M. Liu, S. Lee, J. Hicks, J. Ankorn, M. King, S. Xu, and Arvind, "Bluedbm: An appliance for big data analytics," in *Proceedings of the 42Nd Annual International Symposium on Computer Architecture*, ser. ISCA '15. New York, NY, USA: ACM, 2015, pp. 1–13.
- [7] M. Blott, K. Karras, L. Liu, K. Vissers, J. Bär, and Z. István, "Achieving 10gbps line-rate key-value stores with fpgas," in *Presented as part of the 5th USENIX Workshop on Hot Topics in Cloud Computing*. Berkeley, CA: USENIX, 2013.
- [8] I. T. Association, *Infiniband*, 2014 (Accessed November 18, 2014). [Online]. Available: <http://www.infinibandta.org>
- [9] M. Alizadeh, A. Greenberg, D. A. Maltz, J. Padhye, P. Patel, B. Prabhakar, S. Sengupta, and M. Sridharan, "Data center tcp (dctcp)," in *Proceedings of the ACM SIGCOMM 2010 Conference*, ser. SIGCOMM '10. New York, NY, USA: ACM, 2010, pp. 63–74.
- [10] A. Theodore Markettos, P. Fox, S. Moore, and A. Moore, "Interconnect for commodity fpga clusters: Standardized or customized?" in *Field Programmable Logic and Applications (FPL)*, 2014 24th International Conference on, Sept 2014, pp. 1–8.
- [11] T. Bunker and S. Swanson, "Latency-optimized networks for clustering fpgas," in *Field-Programmable Custom Computing Machines (FCCM)*, 2013 IEEE 21st Annual International Symposium on, April 2013, pp. 129–136.
- [12] K. E. Fleming, M. Adler, M. Pellauer, A. Parashar, Arvind, and J. Emer, "Leveraging latency-insensitivity to ease multiple fpga design," in *Proceedings of the ACM/SIGDA International Symposium on Field Programmable Gate Arrays*, ser. FPGA '12. New York, NY, USA: ACM, 2012, pp. 175–184.

UNCLASSIFIED

UNCLASSIFIED